

*Statistics*

≡ INDEX

1 | Descriptive

2 | Probability

3 | Discrete Random Variables

4 | Continuous Random Variables

5 | Sample Distribution

6 | Estimation

7 | Testing Statistical Hypothesis

8 | About 2 Population

9 |

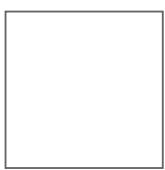
10 | Analysis of Variance

11 | Chi-Squared Goodness of Fit Test

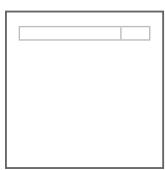
12 | Simple Linear Regression

+  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12

**+** **TEMPLATES** Copy and paste in the section of your choice.



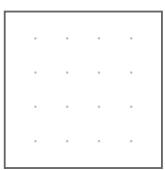
Blank



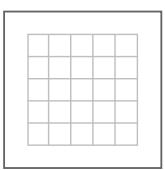
Blank + Title



Lines



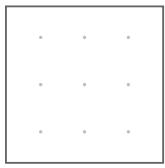
Dots



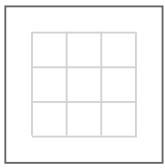
Grid



Wide Lines



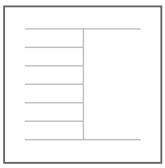
Wide Dots



Wide Grid



Top Blank  
Bottom Lined



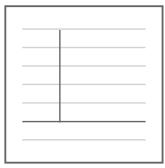
Left Lined  
Right Blank



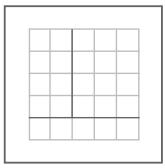
1 Column Lined



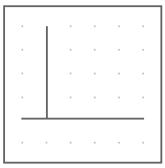
2 Column Lined



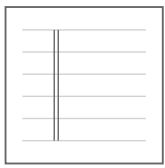
Cornell



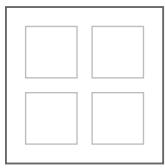
Cornell Grid



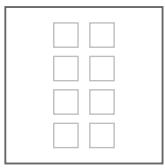
Cornell Dot



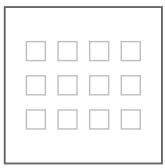
Legal



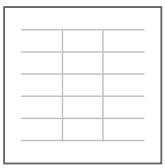
Boxes x 4



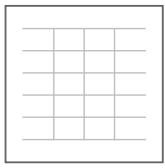
Boxes x 8



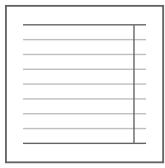
Boxes x 12



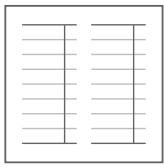
3 Column Table



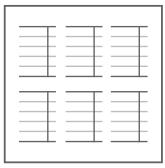
4 Column Table



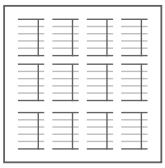
Lists x 1



Lists x 2



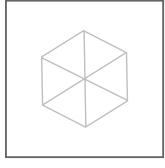
Lists x 6



Lists x 12



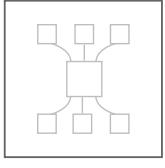
Isometric Grid 1



Isometric Grid 2



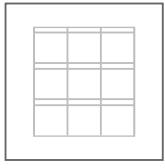
Music Paper



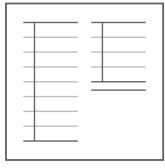
Mind Map



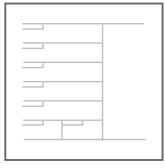
Recipe



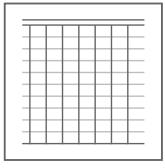
12 Month Plan



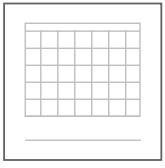
Daily Plan



Weekly Plan



Weekly Plan 2



Monthly Plan

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

...

+

1

2

3

4

5

6

7

8

9

10

11

12

☰ + 1 2 3 4 5 6 7 8 9 10 11 12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

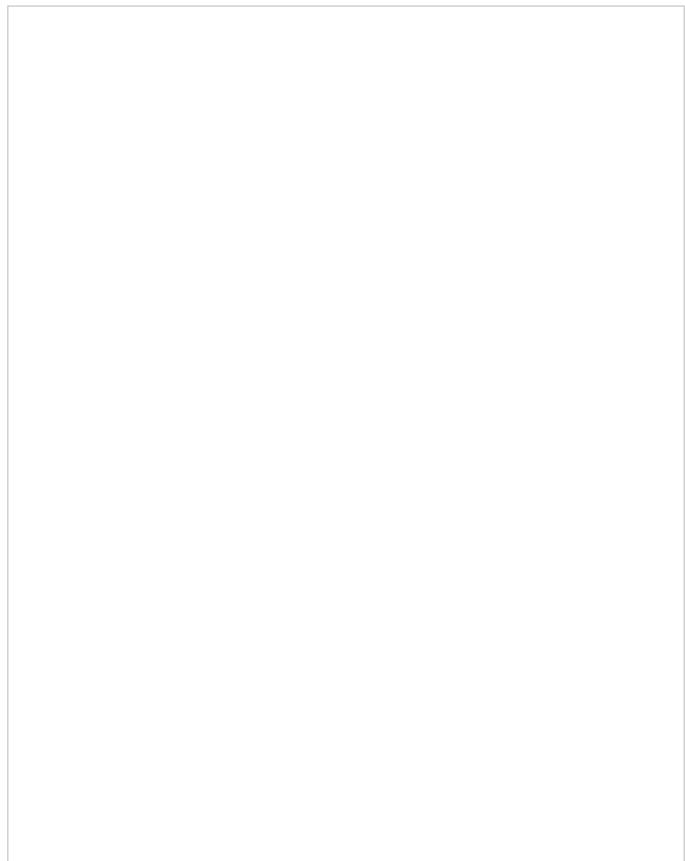
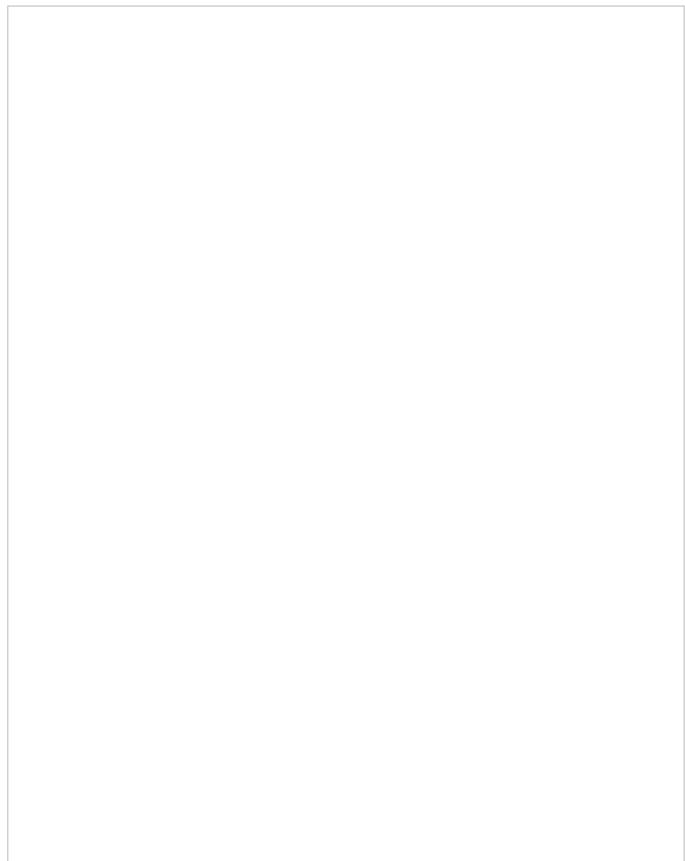
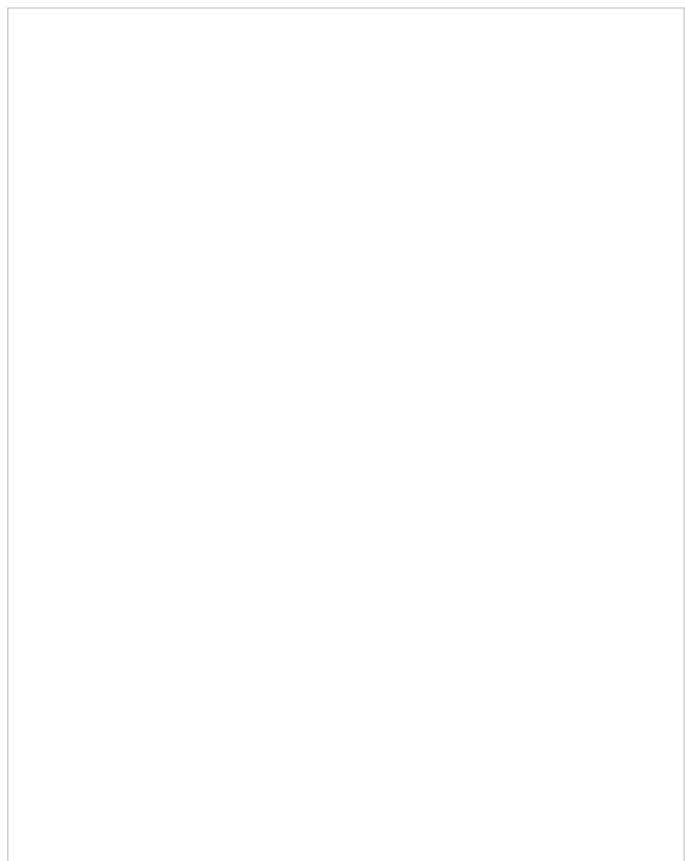
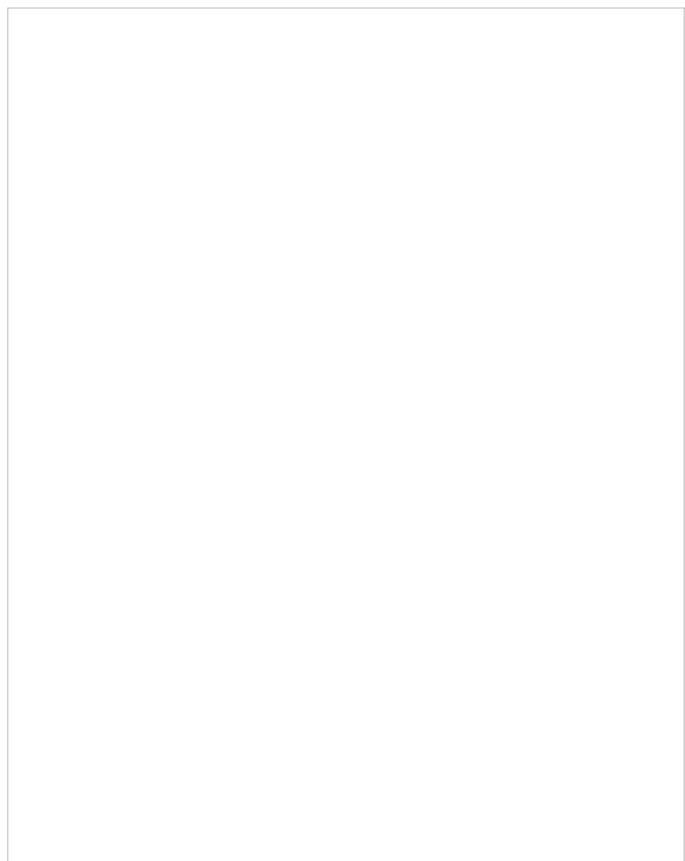
8

9

10

11

12



≡

+

1

2

3

4

5

6

7

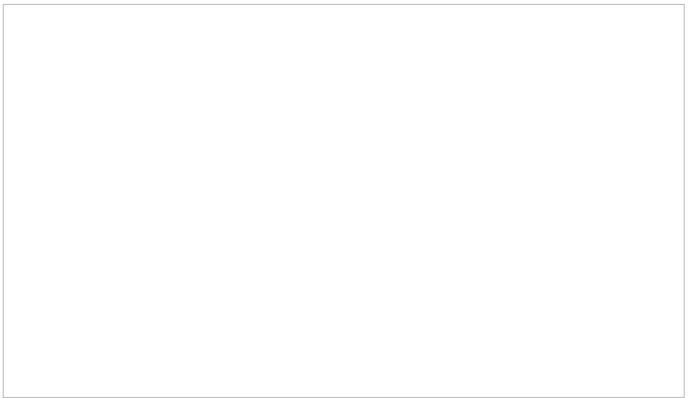
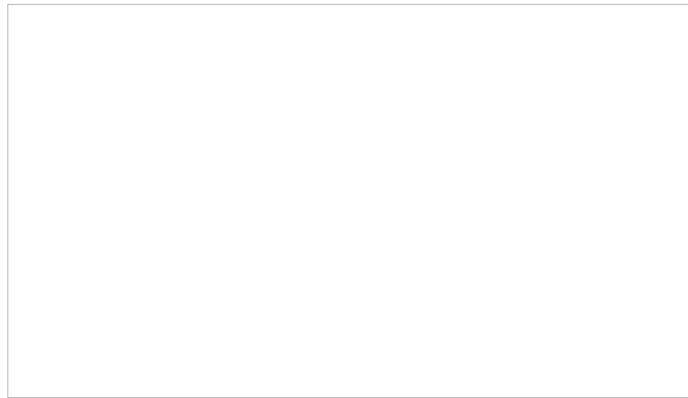
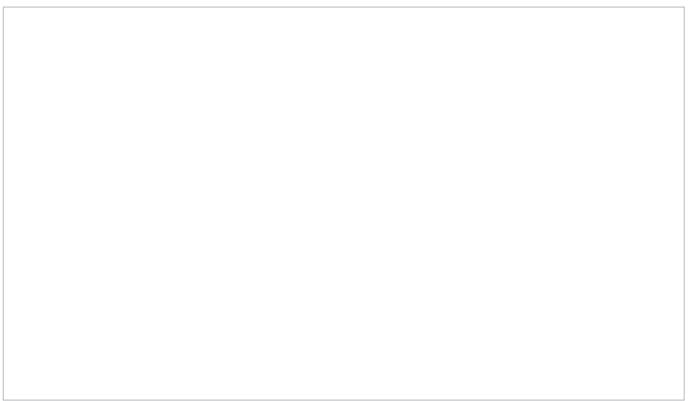
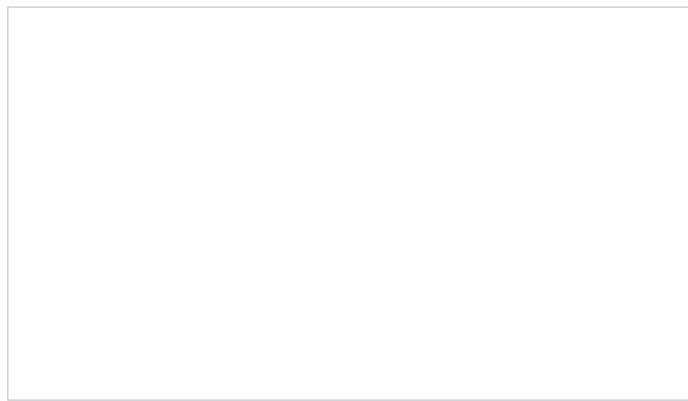
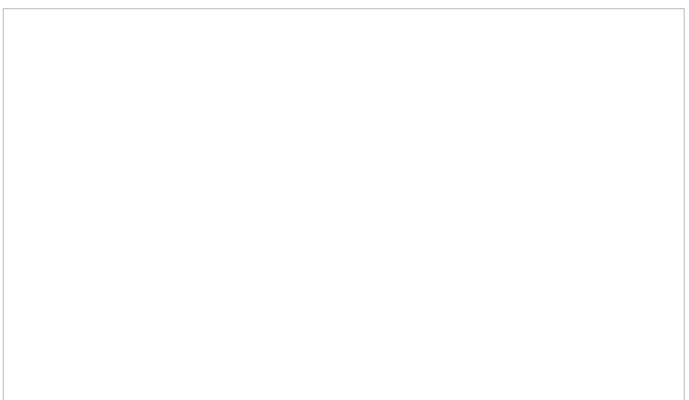
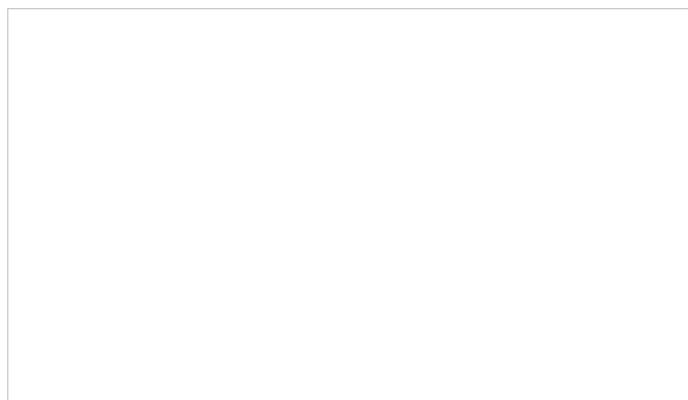
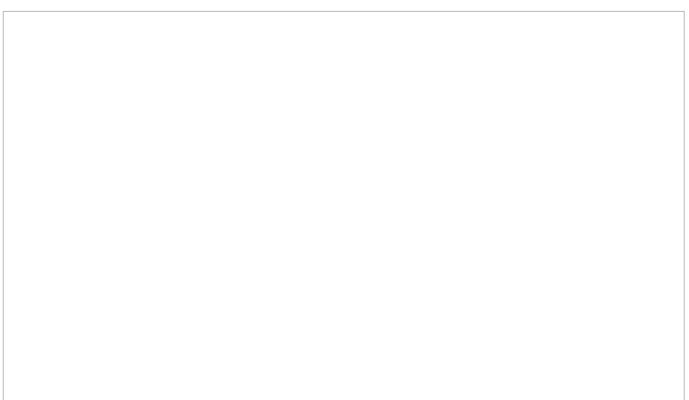
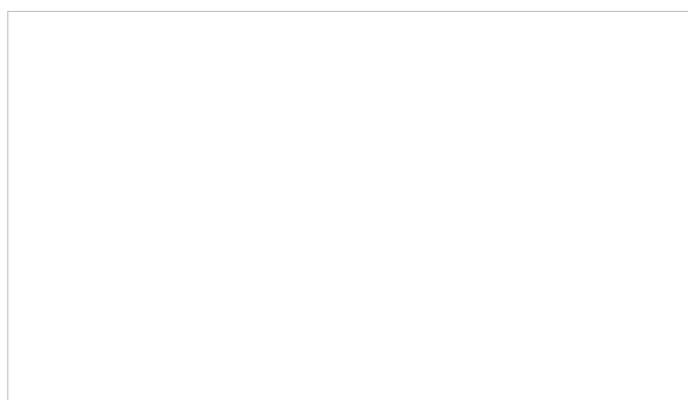
8

9

10

11

12



≡

+

1

2

3

4

5

6

7

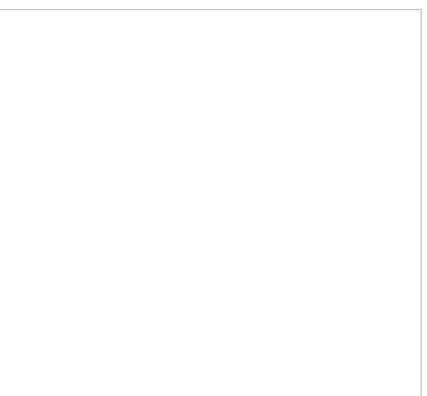
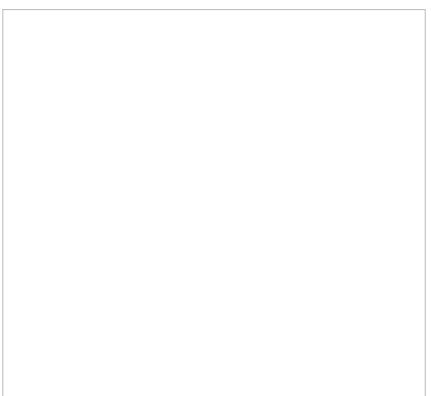
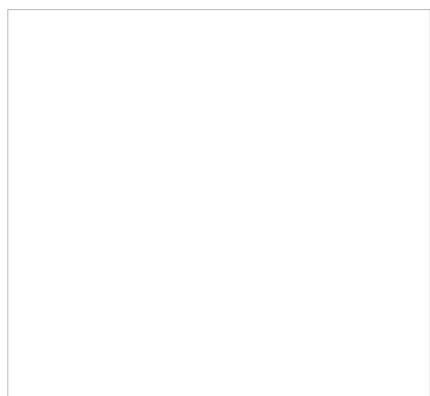
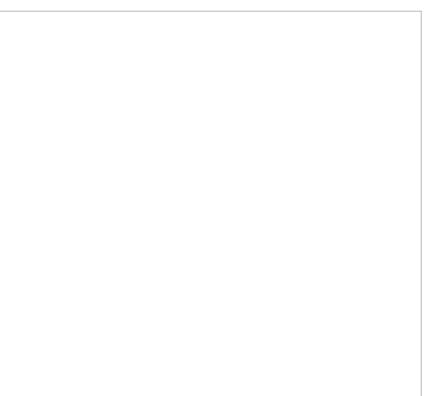
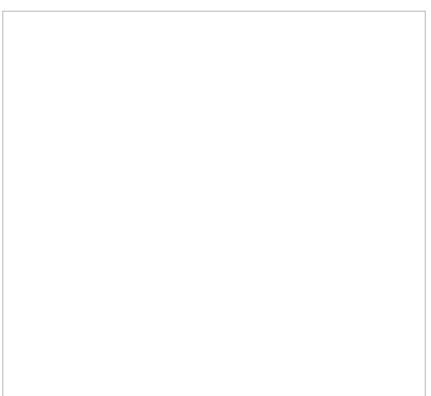
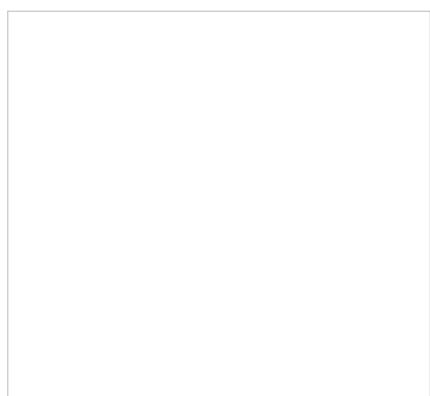
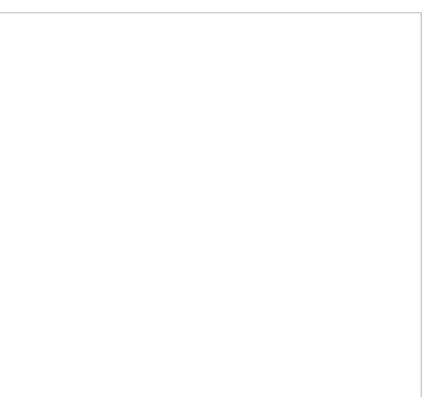
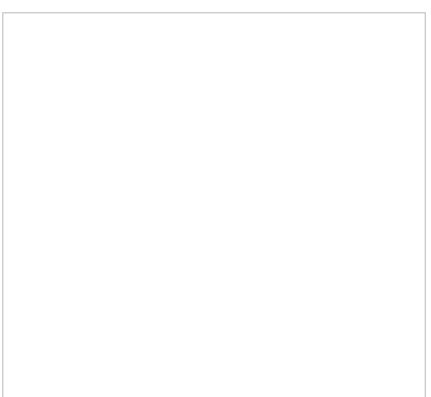
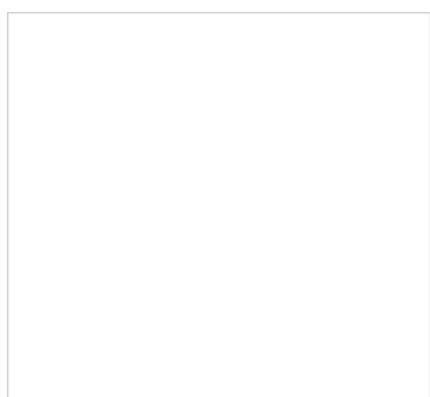
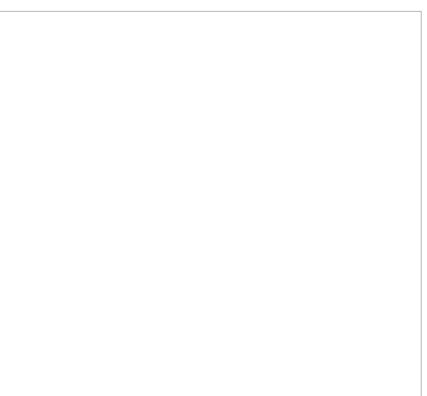
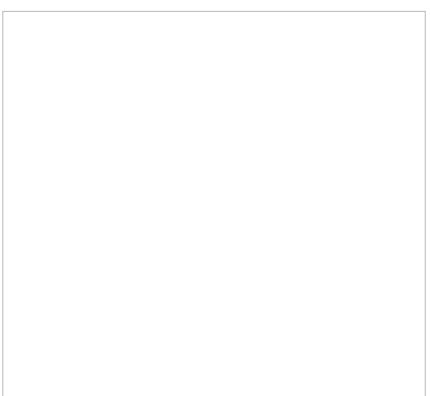
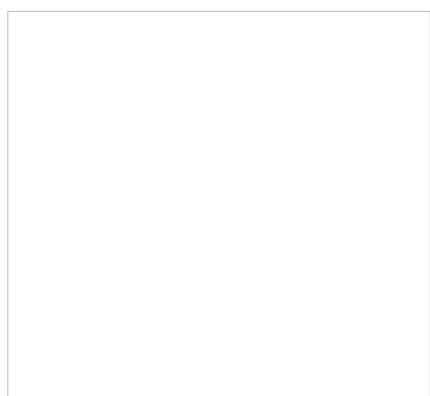
8

9

10

11

12







≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

≡

+

1

2

3

4

5

6

7

8

9

10

11

12

三

+

1

2

3

4

5

6

7

8

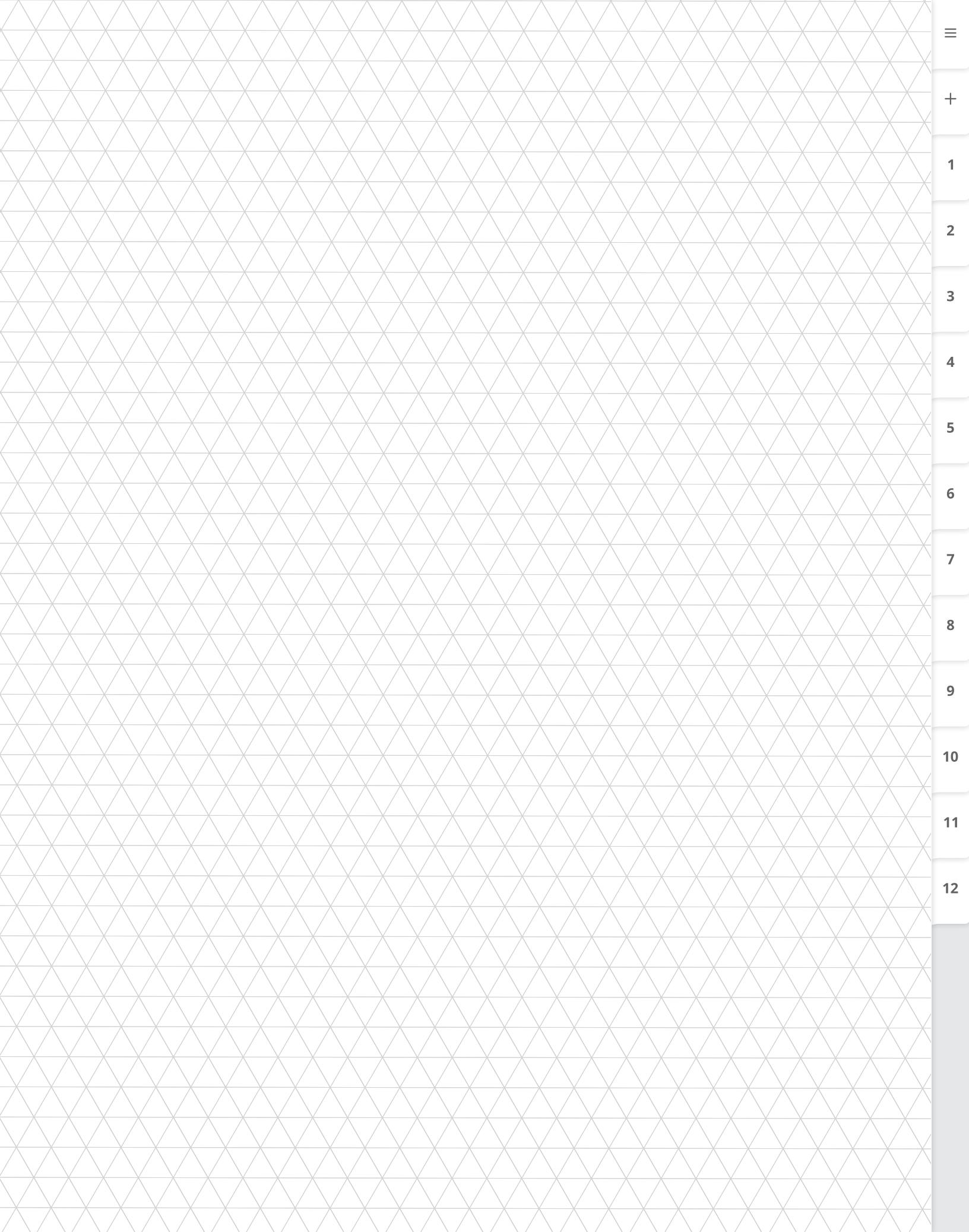
9

10

11

12

☰ + 1 2 3 4 5 6 7 8 9 10 11 12



☰ + 1 2 3 4 5 6 7 8 9 10 11 12

≡

+

1

2

3

4

5

6

7

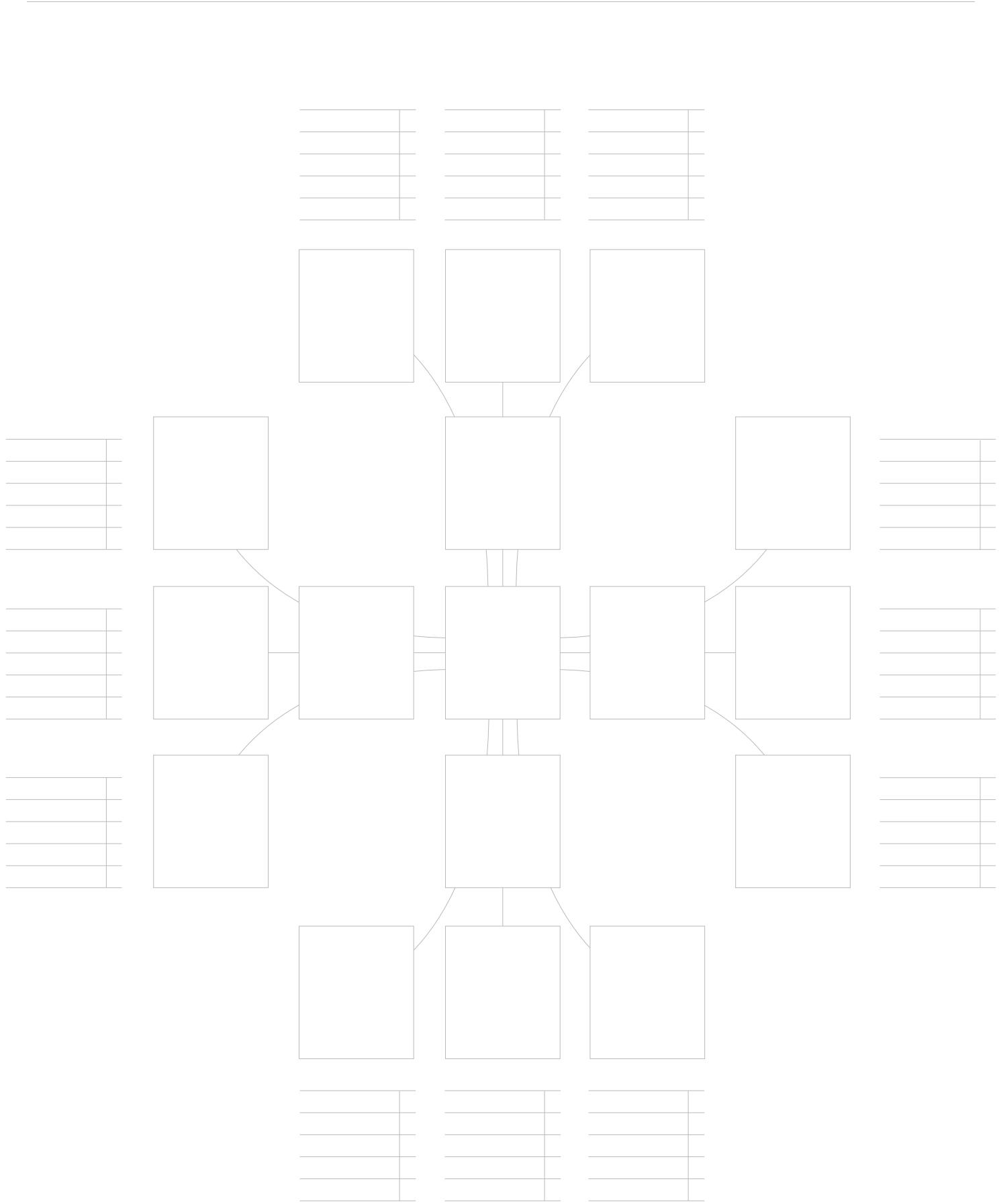
8

9

10

11

12



+

1

2

3

4

5

6

7

8

9

10

11

12

RECIPE FOR

SERVES

PREP TIME

COOK TIME

TOTAL TIME

INGREDIENTS

[Photo Here]

DIRECTIONS

NOTES

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12


1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12



## SCHEDULE

6:00

7:00

8:00

9:00

10:00

11:00

12:00

1:00

2:00

3:00

4:00

5:00

6:00

7:00

8:00

9:00

10:00

11:00

## TASKS

## NOTES

...

+

1

2

3

4

5

6

7

8

9

10

11

12



(⌚)							
6:00							
7:00							
8:00							
9:00							
10:00							
11:00							
12:00							
1:00							
2:00							
3:00							
4:00							
5:00							
6:00							
7:00							
8:00							
9:00							
10:00							
11:00							

+  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12


1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12

1

## Description

# 1. Description

## Measure Location

Median (50%) Mode (most freq.)

$$\text{mid} = Q_2$$

## Percentiles ( $P_i$ )

$$i = |x| \cdot p\% \quad \begin{cases} \text{if } i \in \mathbb{Z}, \quad P_i = \frac{1}{2}(P_i + P_{i+1}) \\ \text{if } i \notin \mathbb{Z}, \quad P_i = P_i \end{cases}$$

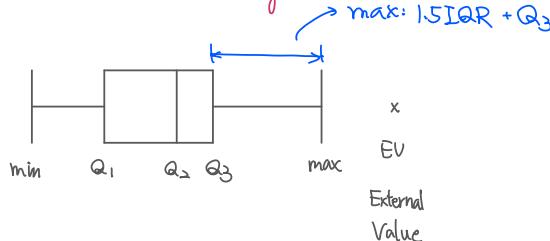
Why?  $p$ th percentiles is mean "How many you win in 100 item"?  
and you wouldn't beat yourself, and the data is decrease.

## Quartiles ( $Q_{1,2,3}$ , IQR)

$$Q_i = P_{25 \cdot i} \quad i=1,2,3$$

$$\text{IQR} = Q_3 - Q_1$$

## 5 number Summary



## Mean ( $\bar{x}$ )

$$\bar{x} = \frac{\sum x_i}{|x|}$$

1. Sum of deviation is "0":  $\sum(x_i - \bar{x}) = 0$

2. Balance point: Center of mass

3.  $\min \arg(\sum(x_i - a)^2) = \bar{x}$

4. If  $y_i = \alpha x_i + \beta \rightarrow \bar{y} = \alpha \bar{x} + \beta$

## Measure Variability

### Sample Variance & Standard Deviation

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$$= \frac{\sum x_i^2 - n\bar{x}^2}{n-1}$$

<tips>

$$\sum \bar{x} = \sum x_i = n\bar{x}$$

$$1. \text{ if } y = \alpha x + \beta \rightarrow S_y^2 = \alpha^2 S_x^2$$

$$\begin{aligned} S^2 &= \frac{1}{n-1} \sum (x_i - \bar{x})^2 \\ &= \frac{1}{n-1} \sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ &= \frac{1}{n-1} \left( \sum x_i^2 - 2 \sum x_i \bar{x} + \sum \bar{x}^2 \right) \\ &= \frac{1}{n-1} \left( \sum x_i^2 - 2 n \bar{x} \bar{x} + n \bar{x}^2 \right) \\ &= \frac{1}{n-1} (\sum x_i^2 - n \bar{x}^2) \end{aligned}$$

### Coefficient Variation

$$CV = \frac{S}{\bar{x}} \cdot 100\%$$

Eliminate the value effect

$$\text{eg. } \bar{x} = 10, S_x = 30 \Rightarrow CV_x = 300\% \quad \text{M}$$

$$\bar{y} = 1000, S_y = 300 \Rightarrow CV_y = 33.3\%$$

### Sample Covariance (with value effect)

$$\text{Cov}_{xy} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y}) \Rightarrow \frac{1}{n-1} \sum x_i y_i - n \bar{x} \bar{y}$$

### Sample Correlation Coefficient (-1 < r < 1)

$$r = \frac{\text{Cov}_{xy}}{S_x S_y}$$

Eliminate the value effect

<Why>

$$\$ \quad \mathbf{x} = [x_1, \dots, x_n]^T$$

$$\text{the Euclidean Norm: } \|x\|_2 = \sqrt{\sum x_i^2}$$

$$\text{and the SD: } \text{SD} = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$$

$$\text{and } \text{cov}_{xy} = \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$$

$$\begin{aligned} \Rightarrow r_{xy} &= \frac{\langle x - \mu_x, y - \mu_y \rangle}{\sqrt{\langle x - \mu_x, x - \mu_x \rangle} \sqrt{\langle y - \mu_y, y - \mu_y \rangle}} \\ &= \frac{\langle x - \mu_x, y - \mu_y \rangle}{\|x - \mu_x\| \|y - \mu_y\|} \end{aligned}$$

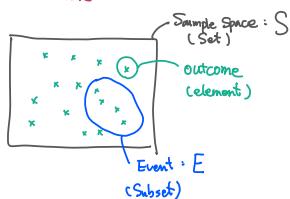
2

# Probability

## 2. Probability

### Intro

#### Some Terms



#### Assigning Prob. to Event

$$P(E) = \frac{\#E}{\#S} \text{ sample Space}$$

#### Rules of Prob.

1. if  $A \in \text{Event in an Exp.} \rightarrow A^c = \{x : x \notin A\}$

A	$A^c$
---	-------

2. if  $AB = \emptyset \rightarrow A, B$  are Mutually Exclusive (Disjoint)

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

#### Conditional Prob.

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \begin{matrix} \text{new event} \\ \text{new Sample Space} \end{matrix}$$

< Dependent & Independent >

if  $P(A|B) = P(A) \rightarrow A, B$  are indep.  
 $\therefore P(A \cap B) = P(A) \cdot P(B)$

#### Multiplication Rule

1.  $P(B|A) = \frac{P(AB)}{P(A)} \Rightarrow P(AB) = P(B|A) \cdot P(A)$

2.  $P(C|BA) = \frac{P(CBA)}{P(BA)} \Rightarrow P(ABC) = P(C|BA) \cdot P(B|A) \cdot P(A)$

2.  $P(DCBA) = P(D|CBA) \cdot P(C|BA) \cdot P(B|A) \cdot P(A)$

#### Bayes' Theorem

if  $A_1 \cup A_2 = S \Rightarrow A_1, A_2 = \emptyset$  ( $A_1 | A_2$ ) (Mut. Exc.)

$$P(A_1 | B) = \frac{P(A_1) \cdot P(B | A_1)}{P(A_1) \cdot P(B | A_1) + P(A_2) \cdot P(B | A_2)}$$

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

< tips >

$P(A_1) \cdot P(A_2)$ : prior prob. (Before)

$P(A_1 | B)$ : posterior prob. (After)

eg.

70% of the proposed zoning change proposals were approved.  
 20% approve, but received negative recommendations  
 10% reject, but received positive recommendation  
 $Q: P(\text{reject})$

step 1)  $P(A_1) \cdot P(A_2)$ : Prior  
 let  $\begin{cases} A_1: \text{approve} \\ A_2: \text{reject} \end{cases}$

$$70\% = P(A_1), 30\% = P(A_2)$$

$$P(A_1) \cdot P(A_2) = 0.21$$

step 2) New Info:  
 let  $B = \text{received "negative" recommendation}$

$$\text{step 3) Bayes' Thm}$$

$$20\% = P(B | A_1)$$

$$10\% = P(B | A_2)$$

step 4) Posterior  $P(A_1 | B), P(A_2 | B)$

$$P(A_1 | B) = \frac{P(A_1) \cdot P(B | A_1)}{P(B)}$$

$$= \frac{P(A_1) \cdot P(B | A_1)}{P(A_1) \cdot P(B | A_1) + P(A_2) \cdot P(B | A_2)}$$

$$= \frac{.7 \cdot .2}{.7 \cdot 2 + .3 \cdot .9}$$

< Tabular Approach >

Event $A_i$	Prior $\times$ Condition	= Joint	Posterior
	$P(A_i)$	$P(B   A_i)$	$P(A_i   B)$
$A_1$	.7	.2	$\frac{P(A_1 \cdot B)}{P(B)}$
$A_2$	.3	.9	$\frac{P(A_2 \cdot B)}{P(B)}$

$$P(B) = P(A_1 \cdot B) + P(A_2 \cdot B)$$

$$P(B) = .41$$

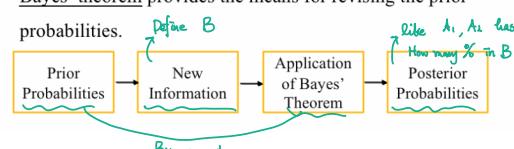


$$P(B) = \underbrace{P(A_1) \cdot P(B | A_1)}_{\text{when } A_1, P(B)} + \underbrace{P(A_2) \cdot P(B | A_2)}_{\text{when } A_2, P(B)}$$

### Bayes' Theorem

$$P(A_1) \cdot P(A_2)$$

- Often we begin probability analysis with initial or prior probabilities.
- Then, from a sample, special report, or a product test we obtain some additional information.
- Given this information, we calculate revised or posterior probabilities.
- Bayes' theorem provides the means for revising the prior probabilities.



By prompt.

Step 2) New Info:  
 let  $B = \text{received "negative" recommendation}$

Step 4) Posterior  $P(A_1 | B), P(A_2 | B)$

$$P(A_1 | B) = \frac{P(A_1) \cdot P(B | A_1)}{P(B)}$$

$$= \frac{P(A_1) \cdot P(B | A_1)}{P(A_1) \cdot P(B | A_1) + P(A_2) \cdot P(B | A_2)}$$

$$= \frac{.7 \cdot .2}{.7 \cdot 2 + .3 \cdot .9}$$

3

## *Discrete Random Variable*

☰  
+  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12

### 3. Discrete Random Variable

#### Random Variable

##### Random Variable

a value determined by the outcome  
of an experiment.

Denoted:  $P(X=x)$ ,  $P$  is Prob. func.

e.g. toss a coin 3 times

let  $X$  be the number of head  
then  $X$  is a RV

$$\left. \begin{array}{l} P(X=0) = \frac{1}{8} \\ P(X=1) = \frac{3}{8} \\ P(X=2) = \frac{3}{8} \\ P(X=3) = \frac{1}{8} \end{array} \right\} \Rightarrow P(X=x) = \frac{1}{8} \binom{3}{x}, x=0,1,2,3$$

#### Probability Distribution

##### Expectation ( $\mu$ )

$$E(X) = EX = \sum x \cdot P(x) \quad \text{Expectation} = \text{Expected Value} = \text{Mean}$$

##### Variance & Standard Deviation ( $\sigma^2, \sigma$ )

$$\sigma^2 = \text{Var}(X) = \sum (x - \mu)^2 \cdot P(x) = EX^2 - \mu^2$$

##### < Some Properties >

1. if  $y = \alpha x + \beta$ ,  $x, y$  are RV

$$\rightarrow EY = \alpha \cdot EX + \beta$$

$$\rightarrow \text{Var} Y = \alpha^2 \cdot \text{Var} X$$

2. standardization:

$$\rightarrow X_{\text{std}} = \frac{X - \mu}{\sigma}$$

3.

$$E(X+y) = EX + EY$$

if  $x, y$  are indep.

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

## Discrete RV

Uniform Dist.  $[X \sim \text{Uni}(1, n)]$

$$P(X) = \begin{cases} \frac{1}{n}, & \text{for } X=1, 2, \dots, n \\ 0, & \text{o.w.} \end{cases}$$

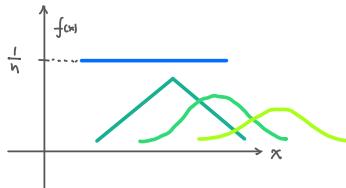
is a const.

<properties>

$$\mu = n \cdot p \Rightarrow \frac{1}{n} \cdot \sum x$$

$$\sigma^2 = n \cdot pq \Rightarrow \frac{1}{n} \cdot \sum (x - \mu)^2$$

$$= \frac{1}{n} (E(x^2) - \mu^2)$$



e.g. Toss a coin,  $P(\text{Head}) = p$

$$\begin{aligned} X &= \# \text{Head in 1 toss} \\ \Rightarrow X &\sim \text{Uni}(1, p) \quad \left\{ \begin{array}{l} \mu_X = \frac{1}{n} \cdot \sum x = 1 \cdot p = p \\ \sigma_X^2 = \frac{1}{n} (E(x^2) - \mu^2) = \frac{1}{n} \left( \frac{1}{2} p^2 + p^2 - p^2 \right) = p - p^2 = p(1-p) \end{array} \right. \end{aligned}$$

e.g. Toss a coin 4 times,  $P(\text{Head}) = p$

$$\begin{aligned} X &= \# \text{Head in 4 tosses} \\ \Rightarrow X &\sim \text{Uni}(4, p) \quad \left\{ \begin{array}{l} \mu_X = 4p \\ \sigma_X^2 = 4pq \end{array} \right. \end{aligned}$$

Binomial Dist  $[X \sim B(n, p)]$

$$P(X) = \binom{n}{x} p^x q^{n-x}, \quad x=0, 1, \dots, n$$

"1 object do M Times"

and only 2 types outcome

and ask occur "How Many Times"

<properties>

$$\mu = n \cdot p$$

$$\sigma^2 = n \cdot pq$$

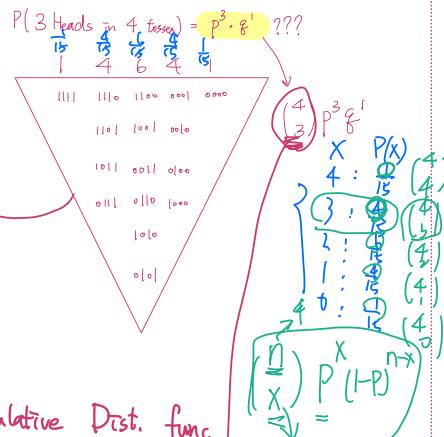
<sup>1</sup> Bernoulli Dist  $[X \sim \text{Ber}(p) = B(1, p)]$

<sup>2</sup> let  $X \sim B(n_x, p), Y \sim B(n_y, p)$

$$\rightarrow X+Y \sim B(n_x+n_y, p)$$

But! if I ask  $P(i \text{ Head in 4 tosses})$

$$P(4 \text{ Head in 4 tosses}) = p^4$$



Cumulative Dist. func.

$$P(X \leq x) = \sum_{i=0}^x P(X_i)$$

$$\begin{aligned} 2: & \quad \begin{array}{c} \text{H} \\ \text{H} \\ \text{T} \\ \text{T} \end{array} / \begin{array}{c} \text{H} \\ \text{H} \\ \text{H} \\ \text{T} \end{array} \\ 3: & \quad \begin{array}{c} \text{H} \\ \text{H} \\ \text{H} \end{array} / \begin{array}{c} \text{H} \\ \text{H} \\ \text{T} \\ \text{T} \end{array} \\ 4: & \quad \begin{array}{c} \text{H} \\ \text{H} \\ \text{H} \\ \text{H} \end{array} / \begin{array}{c} \text{H} \\ \text{H} \\ \text{H} \\ \text{T} \\ \text{T} \end{array} \\ 5: & \quad \begin{array}{c} \text{H} \\ \text{H} \\ \text{H} \\ \text{H} \\ \text{H} \end{array} / \begin{array}{c} \text{H} \\ \text{H} \\ \text{H} \\ \text{H} \\ \text{T} \\ \text{T} \end{array} \end{aligned}$$

$$\begin{array}{c} \binom{2}{0} \binom{2}{1} \binom{2}{2} \\ \binom{3}{0} \binom{3}{1} \binom{3}{2} \binom{3}{3} \\ \binom{4}{0} \binom{4}{1} \binom{4}{2} \binom{4}{3} \binom{4}{4} \\ \binom{5}{0} \binom{5}{1} \binom{5}{2} \binom{5}{3} \binom{5}{4} \binom{5}{5} \end{array}$$

## Hypergeometric Dist. ( $X \sim H()$ )

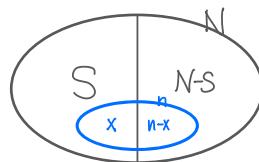
$$P(X) = \frac{\binom{S}{x} \binom{N-S}{n-x}}{\binom{N}{n}}$$

Total

"**2 types** do **N Time**"

and only **2 types** outcome

and ask occur "**How Many Times**"



e.g. 10 red, 20 white balls, draw 5 balls at once.

$$X = \# \text{ red balls}$$

$$P(X) = \frac{\binom{10}{x} \binom{20}{5-x}}{\binom{30}{5}}$$

Total

< properties >

$$\mu = np$$

$$P = \frac{S}{N} \Rightarrow \text{prob. of success}$$

$$\sigma^2 = C \cdot npq$$

$$C = \frac{N-n}{N-1} \Rightarrow C=1 \text{ as } n=1 \text{ or } N \rightarrow \infty$$

correction coeff.

## Geometric Dist. ( $X \sim G(p)$ )

$$P(X) = p q^{X-1}$$

**Until Event Happen**

< properties >

$$\mu = \frac{1}{p}$$

$$\sigma^2 = \frac{q}{p^2}$$

## Poisson Dist. ( $X \sim \text{Poi}(\lambda)$ )

$$P(X) = \lambda^x \frac{e^{-\lambda}}{x!}$$

# Occurrence over a Specified Interval

e.g. ave # accident occurring weekly on a highway is 1.2

Q: P(at least 1 accident this week)

Using  $B(n, p)$

$\Rightarrow X \sim B(n=0, p=1.2)$

of Time or Space

< properties >

$$\mu = \lambda$$

$$\sigma^2 = \lambda$$

1.  $\exists X \sim B(n,p)$ , if  $n \rightarrow \infty$

$\rightarrow X \sim Poi(\lambda)$ ,  $\lambda = np$

2.

very similar with Exponential Dist.

but Poi. Dist. is Discrete

$$X = \sum_{n=1}^{\infty} np$$

Using  $Poi(\lambda)$  to approach

$$\Rightarrow X \sim Poi(\lambda=1)$$

$$X = \lambda^x \frac{e^{-\lambda}}{x!}$$

$$P(X > 0) = P(A) - P(X=0)$$

$$= 1 - \lambda^0 \frac{e^{-\lambda}}{0!}$$

$$= 1 - e^{-\lambda}$$

4

## Continuous Random Variables

## 4. Continuous Random Variables

Continuous RV

Prob. Density Function ( $1 = \int_{-\infty}^{\infty} f(x) dx$ )

$X \sim CRV$

The pdf of  $X$ :

$$P(\alpha < X < \beta) = \int_{\alpha}^{\beta} f(x) dx$$

Expectation

$$EX = \int_{\alpha}^{\beta} x \cdot f(x) dx$$

Variation

$$\text{Var}X = EX^2 - \underline{EX} \mu$$

Uniform Dist. [ $X \sim \text{Uni}(\alpha, \beta)$ ]

$$f(x) = \begin{cases} \frac{1}{\beta - \alpha}, & \alpha < x < \beta \\ 0, & \text{o.w.} \end{cases}$$

< properties >

$$\mu = \frac{\beta + \alpha}{2} \quad (\text{pf 1})$$

$$\sigma^2 = \frac{(\beta - \alpha)^2}{12} \quad (\text{pf 2})$$

Trigonometry Dist. [ $X \sim \text{Tri}( )$ ]

$$\int_{\min}^{\max} f(x) dx = \triangle = 1$$

(pf 1)

$$\begin{aligned} \mu &= \int_{\alpha}^{\beta} x \cdot f(x) dx \\ &= \frac{1}{3} \int_{\alpha}^{\beta} x \cdot \frac{1}{\beta - \alpha} dx \\ &= \frac{1}{3} \left( \frac{\beta^2 - \alpha^2}{\beta - \alpha} \right) \frac{1}{\beta - \alpha} \\ &= \frac{1}{3} (\beta + \alpha) (\beta - \alpha) \frac{1}{\beta - \alpha} \\ &= \frac{1}{3} (\beta + \alpha) * \end{aligned}$$

(pf 2)

$$\begin{aligned} \sigma^2 &= EX^2 - \mu^2 \\ EX^2 &= \int_{\alpha}^{\beta} x^2 \cdot f(x) dx \\ &= \frac{1}{3} \int_{\alpha}^{\beta} x^2 \cdot \frac{1}{\beta - \alpha} dx \\ &= \frac{1}{3} \left( \frac{\beta^3 - \alpha^3}{\beta - \alpha} \right) \frac{1}{\beta - \alpha} \\ &= \frac{1}{3} (\beta^2 + \beta\alpha + \alpha^2) \\ &\therefore \sigma^2 = \frac{1}{3} (\beta^2 + \beta\alpha + \alpha^2) - \frac{1}{3} (\beta + \alpha)^2 \\ &= \frac{1}{3} (\beta + \alpha)^2 - \beta\alpha - \frac{1}{3} (\beta + \alpha) \\ &= \frac{1}{12} (\beta - \alpha)^2 * \end{aligned}$$

Normal Dist. [ $X \sim N(\mu, \sigma^2)$ ]

$$f(x) = \frac{1}{\sqrt{2\pi} \cdot \sqrt{2}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

standardization ( $Z$ -score)

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

<properties>

$$\mu = \mu$$

percentile of  $N(\lambda, \sigma)$

$$\sigma^2 = \sigma^2$$

1.  $Z_n = -Z_{1-n}$ ,  $n \in (0, 1)$

$P(0 < X < n) \approx 0.5 - P(n < X) \approx$

$$Z_{0.95} = 1.96$$

$$Z_{0.05} = 1.645$$

$$Z_{0.1} = 1.28$$

2. let  $X \sim N(\mu_x, \sigma_x^2)$ ,  $Y \sim N(\mu_y, \sigma_y^2)$

(1)  $X + Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$

(2)  $\alpha X + \beta Y \sim N(\alpha \mu_x + \beta \mu_y, \alpha^2 \sigma_x^2 + \beta^2 \sigma_y^2)$

(3)  $X - Y \sim N(\mu_x - \mu_y, \sigma_x^2 + \sigma_y^2)$

Exponential Dist. [ $X \sim \text{Exp}(\lambda)$ ]

Relation to Poisson Dist.

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

if # arrivals in time interval  $(0, t)$

follow a  $\text{Poi}(\lambda \cdot t)$ ,  $\lambda$  is the ave # arrivals per time

$\Rightarrow \begin{cases} N_t : \# \text{ arrival in time interval } (0, t) \\ X : \text{the time, to the 1st arrival} \end{cases}$

$$P(X < t) = P(N_t > 0) = 1 - e^{-\lambda t}$$

$$P(X > t) = P(N_t = 0) = e^{-\lambda t}$$

1. cdf:  $F(x) = P(X \leq x) = 1 - e^{-\lambda x}, x \geq 0$

$$\Rightarrow P(X > t) = e^{-\lambda t}, t > 0$$

2.  $\begin{cases} P_{\text{bin}} : P(\# \text{ events in a time interval}) \\ \text{Exp.} : P(\text{time till 1st event}) \text{ or} \end{cases}$

$P(\text{time between events})$

5

# Sampling Distribution

## 5. Sampling Distribution

### Simple Random Sampling

#### Finite Population

- Sampling with replacement
- Sampling without replacement: most often

#### Infinite Population

- Select from the same popn.
- Each ele. is indep.

#### TID (Independent & Identical Dist.)

- Drawn from the same prob. Dist.
- They constitute a sample from the Dist.

### Dist of the sample mean ( $\bar{X}$ )

#### Finite popu. ( $N < 50n$ )

$$E\bar{X} = \mu$$

$$\text{Var}\bar{X} = \frac{N-n}{N-1} \frac{\sigma^2}{n}, \quad \sigma^2 = \text{Var } X$$

#### Infinite popu. ( $N \geq 50n$ )

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$E\bar{X} = \mu$$

$$\text{Var}\bar{X} = \frac{\sigma^2}{n}$$

## CLT: Central Limit Theorem

### CLT

Let  $\{X_1, \dots, X_n\}$  be a sample from a popu. <sup>(infinite)</sup>

$$\bar{X} \approx N(\mu, \frac{\sigma^2}{n})$$

or

$$\sum X_i \approx N(n\mu, n\sigma^2)$$

for finite popu.

$$\bar{X} \approx N(\mu, C \frac{\sigma^2}{n}), \quad C = \frac{N-n}{N-1}$$

or

$$\sum X_i \approx N(n\mu, C\sigma^2)$$

e.g.  $10^4$  automobile policyholders, if expected yearly claim per policyholder is \$250, SD of \$800.

Q:  $P(\text{Total yearly claim exceeds } \$2.8M)$

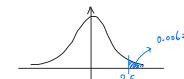
$X_i$ : Yearly claim of policyholder  $i$ ,  $i=1, \dots, 10^4$

$Y = \sum X_i$ : Total yearly claim of all policyholder

$$\stackrel{\text{CLT}}{Y} \approx N(10^4 \cdot \mu, 10^4 \cdot \sigma^2)$$

$$P(Y > 2.8M) = P\left(\frac{Y - \mu_y}{\sigma_y} > \frac{2.8M - \mu_y}{\sigma_y}\right), \text{ where } \begin{cases} \mu_y = 10^4 \mu \\ \sigma_y = \sqrt{10^4} \sigma \end{cases}$$

$$\approx P(Z > 2.5) = 0.0062$$



### Application of CLT

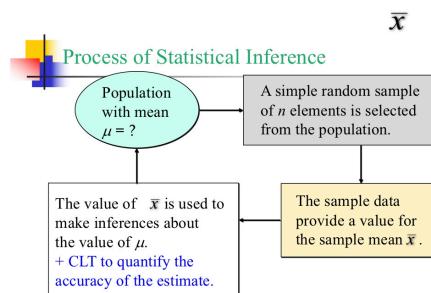
$\bar{X}$  is the sample mean, from popu with  $\mu, \sigma^2$

$$P(\bar{X} < \alpha) = P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq \frac{\alpha - \mu}{\sigma/\sqrt{n}}\right)$$

Standard Normalization

$$\approx P(Z \leq \frac{\alpha - \mu}{\sigma/\sqrt{n}})$$

$$E\bar{X} = \mu, \quad \text{Var}\bar{X} = \frac{\sigma^2}{n}, \quad SD\bar{X} = \frac{\sigma}{\sqrt{n}}$$



## Sampling proportions from a Popn.

### Sample Proportion

- When sampling, there are  $X$  successes from  $Ber(n, p)$ ,  $\hat{p} = \frac{X}{n}$

$$* Ber(p) \Rightarrow B(1, p)$$

$$M = 1 \cdot p, \quad \sigma^2 = 1 \cdot p \cdot q$$

- It's a RV.

Finite vs. Infinite Pop.

When  $n/N$  is small,  $\rightarrow$  Infinite

$X_1, \dots, X_n$  are approximately independent  $Ber(p)$

$\bar{X}$  is often denoted by  $\hat{P}$

When  $n/N$  is large,  $\rightarrow$  Finite

$X_1, \dots, X_n$  are not independent

<properties>

(1) When  $n/N$  is Large,

$$E\hat{P} = p, \text{Var}\hat{P} \approx C \cdot \frac{1}{n} p q$$

$$\frac{\hat{P} - p}{\sqrt{\frac{1}{n} p q}} \approx N(0, 1)$$

(2) When  $n/N$  is Small,

$$E\hat{P} = p, \text{Var}\hat{P} = \frac{1}{n} p q$$

$$\frac{\hat{P} - p}{\sqrt{\frac{1}{n} p q}} \approx N(0, 1)$$

(3) If  $Y \sim B(n, p)$

$\rightarrow Y = X_1 + \dots + X_n$ , where  $X_i$  are i.i.d.  $Ber(p)$  RVs

CLT

$\rightarrow Y \approx N(np, npq)$

$$P(Y \leq k) = P\left(\frac{Y-np}{\sqrt{npq}} \leq \frac{k-np}{\sqrt{npq}}\right) \xrightarrow{\mu} \mu$$

$$\xrightarrow{\sigma} P\left(N(0, 1) \leq \frac{k-np}{\sqrt{npq}}\right)$$

\*  $X \sim B(n, p)$ ,  $\mu = np$ ,  $\sigma^2 = npq$

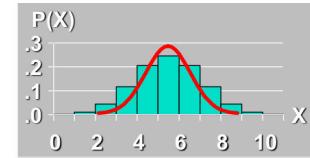
Let  $Y \sim N(\mu, \sigma^2)$

$$P(X \leq x) \approx P(Y \leq x + 0.5)$$

## Normal Approximation of Binomial Distribution

1. Not all binomial tables exist
2. Use normal dist. to approximate
3. Requires large sample size
4. Need correction for continuity

$$n = 10, p = 0.50$$



$x$	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.0000	0.0040	0.0080	0.0110	0.0140	0.0170	0.0199	0.0228	0.0257	0.0279
0.1	0.0000	0.0040	0.0079	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.2	0.0002	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.3	0.0005	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.4	0.0014	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.5	0.0015	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.6	0.0020	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.7	0.0020	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.8	0.0020	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
0.9	0.0015	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.0	0.0012	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.1	0.0008	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.2	0.0005	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.3	0.0003	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.4	0.0002	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.5	0.0001	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.6	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.7	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.8	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
1.9	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.0	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.1	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.2	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.3	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.4	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.5	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.6	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.7	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.8	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
2.9	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352
3.0	0.0000	0.0040	0.0080	0.0117	0.0157	0.0197	0.0236	0.0275	0.0313	0.0352

Comparison

$$X \sim B(5, \frac{1}{2})$$

$$P(X=1) = \binom{5}{1} \frac{1}{2} \left(\frac{1}{2}\right)^4 = 0.156$$

$$P(X \leq \frac{5}{2}) = P(X \leq 2) = \sum_{i=0}^2 \binom{5}{i} \left(\frac{1}{2}\right)^5 = 0.5$$

$$Y \sim N(\mu_x, \sigma_x^2) \Rightarrow \begin{cases} \mu_x = \frac{5}{2} \\ \sigma_x = \sqrt{\frac{5}{4}} \end{cases}$$

$$P(X=1) \approx P(1.05 \leq Y \leq 1.95)$$

$$P(0.5 \leq Y \leq 1.5) = P\left(\frac{0.5 - \mu}{\sigma} \leq Z \leq \frac{1.5 - \mu}{\sigma}\right)$$

$$= P(-1.79 \leq Z \leq -0.9)$$

$$= P(-0.9 \leq Z \leq 1.79)$$

$$= 0.4633 - 0.3159 = 0.15$$

$$P(X=\alpha) \approx P(\underbrace{\alpha-0.5}_{\text{sample size}} \leq Y \leq \underbrace{\alpha+0.5}_{\text{sample size}})$$

$$P(X \leq \frac{\alpha}{2}) = P(X \leq 2) = P(Y \leq \frac{\alpha}{2}) \text{ True}$$

### Chi-Square Dist.

If  $Z_i \sim N(\mu_i, \sigma_i^2)$  are indep.  $i=1, \dots, n$

$$\rightarrow Y = \sum(Z_i)^2 \sim \chi_n^2 \text{ "sample size - 1" = df. "How many N-Dist. sum up."}$$

< properties >

$$EY = n$$

$$\text{Var}Z = 1$$

$$(1) \quad \chi_n^2 \Rightarrow \frac{\sum(X_i - \mu)^2}{\sigma^2}$$

$$(2) \quad \chi_1^2 \Rightarrow \frac{n(X_i - \mu)^2}{\sigma^2} = \left( \frac{X_i - \mu}{\sigma/\sqrt{n}} \right)^2$$

$$(3) \quad \chi_n^2 = \chi_{n-1}^2 + \chi_1^2$$

$$\sum(X_i - \mu)^2 = \sum(X_i - \bar{x})^2 + n(\bar{x} - \mu)^2$$

$$(4) \quad \text{If } X \sim \chi_n^2, Y \sim \chi_m^2$$

$$\rightarrow X + Y \sim \chi_{n+m}^2$$

(5)

let  $X_1, \dots, X_n$  be sample from N-Dist.

$$\text{popu. } \frac{(n-1)S^2}{\sigma^2} = \frac{\sum(X_i - \bar{x})^2}{\sigma^2} \text{ Sample}$$

has a Chi-Square Dist. with df. : (n-1)

(6)

$$\text{Consider } Z_i = \frac{X_i - \mu}{\sigma}, i=1, \dots, n$$

$$\rightarrow \sum\left(\frac{X_i - \mu}{\sigma}\right)^2 = \frac{\sum(X_i - \mu)^2}{\sigma^2} \sim \chi_n^2$$

(7)

If  $X_1, \dots, X_n$  from  $N(\mu, \sigma^2)$  popu.

$$\rightarrow E(S^2) = \sigma^2, S^2 \text{ is sample Variance}$$

6

# Estimation

## 6. Estimation

### Conclusion

Goal: Get "Popu Mean"

How: By "Sample Mean( $\bar{X}$ )" to "Estimate"

Question: "How much" accuracy  $\bar{X}$  is  $\mu_{\text{popu}}$

Tools: CLT

Transform the thinking: By "Interval"

We could say, you give me a  $k\%$ ,

I tell you the interval that Mean would (probability) in it with  $k\%$  confidence.

$$k\% = 1 - \alpha$$

### Process

#### Sample

#### Statistic

proportion

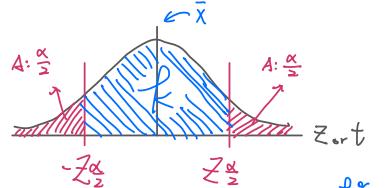
$$\bar{X} \rightarrow np \rightarrow \hat{P}$$

$$\sigma_{\bar{X}}^2 \rightarrow npq \rightarrow n \hat{P} \hat{Q}$$

Sampling  
 $\bar{X}, S^2$

Known  
 $\sigma_{\text{popu}}$ ?

$$\begin{aligned} T &\rightarrow \bar{Z}_{\frac{\alpha}{2}} = \frac{\bar{X} - \mu}{\sigma_{\text{popu}} / \sqrt{n}} \\ F &\rightarrow t_{\frac{\alpha}{2}(n-1)} = \frac{\bar{X} - \mu}{\sigma_{\bar{X}} / \sqrt{n}} \end{aligned}$$



$$P(-Z_{\frac{\alpha}{2}} < SS < Z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$P(-Z_{\frac{\alpha}{2}} < \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} < Z_{\frac{\alpha}{2}})$$

$$P(-Z_{\frac{\alpha}{2}} \cdot \sigma / \sqrt{n} < \bar{X} - \mu < Z_{\frac{\alpha}{2}} \cdot \sigma / \sqrt{n})$$

$$P(\mu \text{ in } \bar{X} \pm Z_{\frac{\alpha}{2}} \cdot \sigma / \sqrt{n}) = k\%$$

<eg>

Weights of RS of 8 people with,  $\sigma_{\text{popu}}: 22$

121, 163, 144, 152,  
183, 130, 128, 160

use the data to estimate popn. 95%

$$1. \text{ SS: } \begin{cases} \bar{X} = 145.5 \\ S = 21.27 \end{cases}$$

$$2. \text{ Known } \sigma_{\text{popu}}, \frac{\alpha}{2} = 0.025, Z_{\frac{\alpha}{2}} = 1.96$$

$$3. [CI: \bar{X} \pm Z_{\frac{\alpha}{2}} \cdot \sigma / \sqrt{n}]$$

$$CI = 145.5 \pm 1.96 \cdot 22 / \sqrt{8}$$

2. If Unknown  $\sigma_{\text{popu}}$

$$3. CI = \bar{X} \pm t_{\frac{\alpha}{2}(n-1)} \cdot \sigma / \sqrt{n}$$

<eg> RS #100, 64%

$$a) \text{ Entire popu. is Event: } E \bar{X} = np = 64 \quad \begin{matrix} \bar{X} = \bar{P} \\ S \end{matrix}$$

b) Standard Error:

[Standard Err = Standard Deviation]

$$S = \sqrt{\frac{npq}{n}} = \sqrt{\frac{0.44 \cdot 0.36}{100}} = 0.098$$

<eg>  $S=11.3, n=81, \bar{X}=74.6, 90\% \text{ CI, popu} \sim N$

$$(1) \bar{X}=74.6, S=11.3, n=81 \quad 0.10 Z_{0.10}$$

$$(2) \text{ popu} \sim N\text{-Dist, } \alpha=0.1, Z_{\frac{\alpha}{2}} = 1.645$$

$$(3) CI = \bar{X} \pm Z_{\frac{\alpha}{2}} \cdot S / \sqrt{n}$$

$$= 74.6 \pm 1.645 \cdot 11.3 / \sqrt{81}$$

$$= (72.53, 76.67)$$

## Point Estimator of Popu. Mean & Variance

### Unbiased

$$E(\bar{X} = p) = p$$

### Popu. "Mean."

$X_1, \dots, X_n$  from N-Dist, and  $\mu$  is unknown.

$$\rightarrow E\bar{X} = \mu_{\text{popu.}}$$

$$P(|\bar{X} - \mu| < 2\sigma/\sqrt{n}) \approx 0.95$$

### <properties>

A bag has " $N$ " balls, " $M$ " are red, a SRS of size  $n$ , let  $X = \# \text{ red}$ .

$$\text{point estimator: } \hat{p} = \frac{X}{n}$$

$$B(n, p)$$

also called "stand. err."  $E\hat{p} = p$ ,  $\text{Var}\hat{p} \approx \frac{pq}{n}$ ,  $p = \frac{M}{N}$

$$SD(\hat{p}) = \sqrt{\frac{pq}{n}}, \text{ if } pq \leq \frac{1}{4}, SD(\hat{p}) \approx \frac{1}{\sqrt{4n}}$$

### Popu. "Variance"

If  $\sigma_{\text{popu.}}^2$  is Unknown,  $X$  sampled from popu.

$$S_{\text{samp}}^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$$

If Known:  $\mu_{\text{popu.}}$

$$\sigma_{\text{popu.}}^2 = \frac{\sum (X_i - \mu)^2}{n}$$

e.g. SRS  $n=50$ , 20 were in favor of proposal

(a) proportion of all students who are in favor

$$E\hat{p} = p = \frac{20}{50} = 0.4 *$$

(b) Stand. err. of this estimate.

$$[SD(\hat{p}) = \frac{pq}{n}]$$

$$SD(\hat{p}) = \sqrt{\frac{0.4 \cdot 0.6}{50}} = 0.07 *$$

now! How much "Confidence" you have that  $SD(\hat{p})$

$$P(|\hat{p} - p| < 2\sqrt{\frac{pq}{n}}) \approx 0.95$$

$$\approx P(|\hat{p} - p| < 2\sqrt{\frac{pq}{n}})$$

SO that  $|p - 0.4| < 2 \cdot 0.07$  with conf. 95%

## Interval Estimators of Mean & Variance

Popu. "Mean"

Interval Estimator :

an interval that is predicted to contain the params.

Margin of Err:

Interval Est. = Point Est.  $\pm$  Margin of Err.

$$M_{\text{popu}} = \bar{x} \pm \text{Margin of Err}$$

## Confidence Interval (CI)

An interval est. states range within

a popu. params. prob. the specified prob.

is called (level of Conf. / Conf. Coef.)  $Z(\cdot)$

Def:  $Z \sim N(0,1) \exists z_\alpha \text{ s.t. } P(Z \geq z_\alpha) = \underline{\alpha}$

- $\alpha = 0.1 \Rightarrow 90\% \text{ CI} \Rightarrow Z_{0.1} = 1.645$
- $\alpha = 0.05 \Rightarrow 95\% \text{ CI} \Rightarrow Z_{0.05} = 1.96$
- $\alpha = 0.01 \Rightarrow 99\% \text{ CI} \Rightarrow Z_{0.01} = 2.58$

## CI of $\mu$

a size  $n$  SRS from inf. prop.  $N(\mu, \sigma^2)$

$\rightarrow 1-\alpha \text{ CI of } \mu: \bar{x} \pm Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$  CI

## <properties>

Somethin we want make  $CI \leq b$

$b$  is a specific value for limit CI.

$$\# \text{ sample: } n \geq \left( \frac{2Z_{\frac{\alpha}{2}}\sigma}{b} \right)^2$$

$$\text{CI: } \bar{x} \pm Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

<pf> When  $n$  is large, by CLT

$$\bar{x} \approx N(\mu, \frac{\sigma^2}{n})$$

$$1-\alpha = P(-Z_{\frac{\alpha}{2}} \leq \frac{\bar{x}-\mu}{\sigma/\sqrt{n}} \leq Z_{\frac{\alpha}{2}})$$

$$= P(\bar{x} - Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}})$$

<eg> sample B  $\sim N(\mu, \sigma^2)$  from A

sample 10 times: 17, 21, 20, 18, 19, 23, 20, 21, 16, 19

95% CI of  $\mu$ ?

$$\bar{x} = \frac{\sum x_i}{n} = 19.3$$

$$\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

$$\Rightarrow 19.3 \pm 1.96 \cdot 3/\sqrt{10} \Rightarrow 19.3 \pm 1.86$$

if want make CI to 0.78

then need to increase How Much # sample

$$\text{CI} = Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

$$0.78 = 1.96 \cdot \frac{3}{\sqrt{n}} \Rightarrow n = 56.8$$

<eg>  $S=0.3$ , wanted 90% certain that the est.

of mean weight is correct within  $\pm 0.1$  pound,  
How Much # sample?

$$[\text{CI: } \bar{x} \pm Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}]$$

$$90\% \Rightarrow \alpha = 0.1 \Rightarrow Z_{0.1} = 1.645$$

$$0.1 = Z_{0.1} \cdot \frac{\sigma}{\sqrt{n}} = 1.645 \cdot \frac{0.3}{\sqrt{n}}$$

$$\therefore n = 24.35$$

## Small Sample Case

Popu is not N-Dist ?

Increase # sample + CLT ( $\bar{x}$ )

Popu  $\sim N(\text{?}, \sigma^2)$

Interval-Est. + CI

Popu  $\sim N(\text{?}, \text{?})$

Student t-Dist + CI

## Student's t-Dist

When # sample is small

$$\frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{n-1}$$

t-Dist has  $\mu=0$ ,  
but s.d. is depend on the # sample

<eg> # sample = 10,  $\mu = 50$ ,  $\sigma = 60$

95% CI est.?

$$[ 95\% \text{ CI} \Rightarrow \alpha = 0.05 ]$$

$$[ \text{CI} : \bar{x} \pm t_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} ]$$

$df = (n-1)$

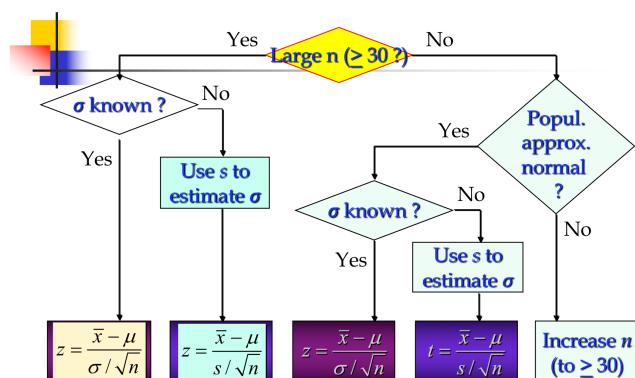
$$550 \pm 2.262 \cdot \frac{60}{\sqrt{10}}$$

## Degree of Freedom (d.f.)

d.f. refers to the # indep. observations

$S_{\text{samp}}^2$  has " $n-1$ " d.f.

## Summary of Interval Estimation Procedures for a Population Mean



## CI of P (popu. proportion)

## CI of P

# sample from a popu. with prop. P  
Suppose Sample proportion is  $\hat{P}$

$$\text{if } N \text{ is large} \rightarrow X \sim B(n, p)$$

$$\begin{array}{c} \text{if } N \text{ is large} \\ \rightarrow X \sim B(n, p) \end{array}$$

$$\text{Var } X = npq$$

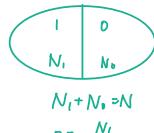
$$\text{Var } \frac{X}{n} = \frac{pq}{n}$$

$$\hat{P} \approx N(p, \frac{pq}{n})$$

$$\text{eg: } \mu = n\hat{p} \Rightarrow \bar{X} = \frac{\mu}{n} = \hat{P}$$

$$\sigma_{\text{pop}}^2 = npq \Rightarrow S = \sqrt{\frac{npq}{n^2}} = \frac{\sqrt{pq}}{n}$$

popu. proportion



$$N_1 + N_0 = N$$

$$p = \frac{N_1}{N}$$

$$\text{Var } X = npq$$

$$\text{Var } \frac{X}{n} = \frac{pq}{n}$$

eg: 220 registered voters, out of 300 contacted  
95% CI Est. for the proportion?

$$\hat{p} = \frac{220}{500} \quad [ \hat{p} \approx N(p, \frac{pq}{n}) ]$$

$$\hat{p} \pm 1.96 \sqrt{\frac{\hat{p}\hat{q}}{n}} \quad [ \hat{p} \pm Z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}\hat{q}}{n}} ]$$

0.975  
0.025

$$0.44 \pm 1.96 \sqrt{\frac{0.44 \cdot 0.56}{500}} \quad [ 95\% \text{ CI} \Rightarrow \alpha=0.05 \Rightarrow Z_{\alpha/2} = 1.96 ]$$

eg: H<sub>0</sub>: p ≤ 0.5  
over  $\frac{1}{2}$ , n=920, M<sub>0.5</sub>=0.52

$$\begin{cases} H_0: p \leq 0.5 \\ H_1: p > 0.5 \end{cases}$$

7

## 7 | Testing Statistical Hypothesis

## F. Testing Statistical Hypothesis

### Conclusion

Goal : accept/reject assumption

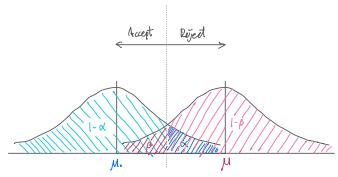
How : By Sampling

Question : When we accept/reject

Tools : CLT,  $\alpha$ ,  $\beta$ , power

→ significant level

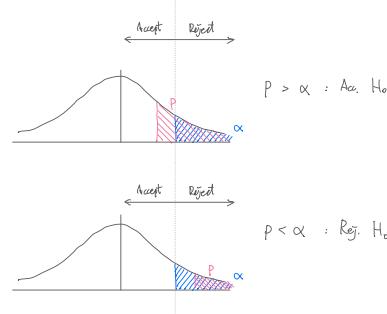
$$\text{Type I error } (\alpha) = P(\text{Rej.} \mid \text{correct})$$



$$\text{Type II error } (\beta) = P(\text{Acc.} \mid \text{incorrect})$$

$$\text{Power} = 1 - \beta$$

$$p\text{-value} = \bar{x} \text{ lead position}$$



### Process

Hypo  
 $H_0, H_1, \mu_0$

Select  
 $\alpha$

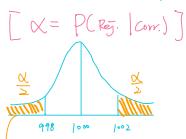
Test Statistic  
 $TS = \frac{\bar{x} - \mu_0}{S/\sqrt{n}}$

Formula Decision  
Critical Region

Make Decision  
Accept / Reject

$$\mu_T = 1000 \rightarrow \text{True value}$$

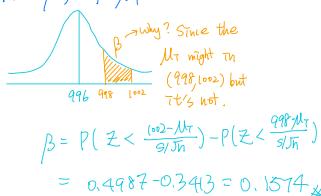
$\left\langle \text{eg} \right\rangle \text{Find } \alpha, H_0 = \mu \in (998, 1002)$   
 $SS: n=36, \bar{x}=998, S=2$



why?  $\mu_0$  might be any value  
if in these area, then we do the type I error!

$$\begin{aligned} \alpha &= 1 - P\left(\left|\frac{\bar{x} - \mu}{S/\sqrt{n}}\right| < \frac{|1002 - 1000|}{S/\sqrt{n}}\right)^{***} \\ &= 1 - P(|Z| \leq 1) \\ &= 2(1 - P(Z < 1)) \\ &= 2(1 - 0.8413) \\ &= 0.3174 \end{aligned}$$

$\left\langle \text{eg} \right\rangle \text{Find } \beta, \text{ if } \mu_1 = 996$



$\left\langle \text{eg} \right\rangle H_0 = 250, \text{Sig. lev.} = 0.05$

$$TS: n=81, \bar{x}=255, S=3$$

$$\begin{cases} H_0: \mu = 250 = \mu_0 \\ H_1: \mu \neq 250 \neq \mu_0 \end{cases}$$

$$\text{step2)} \quad \alpha = 0.05, Z_{\alpha/2} = 1.96$$

$$\text{step3)} \quad TS = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} = \frac{255 - 250}{3/\sqrt{81}} = 15$$

$$15 > 1.96 \Rightarrow \text{Reject!}$$

$\left\langle \text{eg} \right\rangle H_0: \mu \leq 140, \text{Sig. level} = 0.05$

$$TS: n=900, \bar{x}=146, S=48.2$$

$$\begin{cases} H_0: \mu = 140 = \mu_0 \\ H_1: \mu > 140 \neq \mu_0 \end{cases}$$

$$\text{step2)} \quad \alpha = 0.05, Z_{\alpha} = 1.645$$



$$\text{step3)} \quad Z = \frac{146 - 140}{48.2/\sqrt{900}} = 3.726$$

$$\text{step4)} \quad Z = 3.726 > 1.645 = Z_{\alpha}$$

## Statistical Hypothesis

### Hypothesis

a statement about the nature of a popu.

### Hypo. Testing

By Sample Evidence to determine whether the hypo.'s approvability.

## Null & Alternative Hypo.

- null hypo.  $H_0$  is a

Tentative Assumption about a popu. param.

- alt. hypo  $H_1$  is a

Opposite of what's stated in the  $H_0$

Assume: the popu. mean age = 45 ( $H_0$ )

popu.  $\rightarrow$  samp.: check  $M_{\text{samp}}$

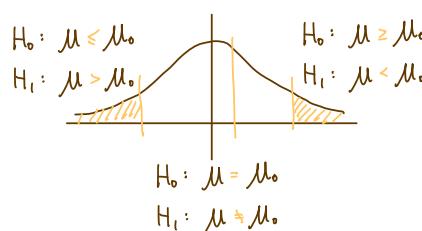
if  $M_{\text{samp}} \approx 45$   
the assume is correct!

But, How is " $\approx$ ",  
And How about the sample bias?

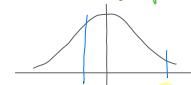
### < properties >

(1) The " $=$ " always appears in  $H_0$

(2) 3 main form



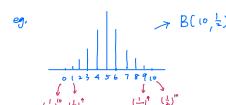
### Reason for Rejecting $H_0$



if  $M_{\bar{x}} = \alpha \rightarrow H_0 \text{ is correct}$

if  $M_{\bar{x}} = \beta \rightarrow H_0 \text{ is incorrect}$

The more  $M_{\bar{x}}$  close to  $H_0$ ,  
the more prob.  $H_0$  is correct



# Type I & Type II Err.

## Type I Err.

Reject  $H_0$ , when it's True

## Type II Err.

Accept  $H_0$ , when it's False

## $\alpha$ -Value

Level of significance :

max allowable prob. of making Type I err.

$$\bullet P(\text{rej. } H_0 \mid H_0 \text{ is True}) \leq \alpha$$

## P-Value

the Smallest Significance Level at which the data lead to Rej  $H_0$

- $P < \alpha \rightarrow \text{Rej. } H_0$
- $P > \alpha \rightarrow \text{Acc. } H_0$

$$\text{P-value} = P(Z \geq Z_0)$$

<eg> previous at least 1.5 mg current less than 1.5 mg

# sample = 20,  $SD = 0.7$ , ave = 1.42

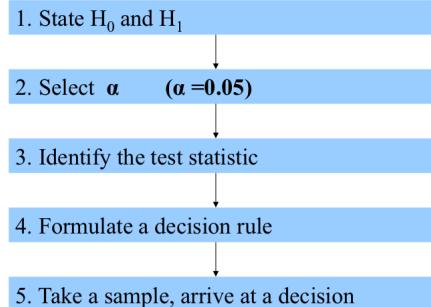
5% level of significance?  $\alpha$

$$1. \begin{cases} H_0: \mu \geq 1.5 \\ H_1: \mu < 1.5 \end{cases}$$

$$2. Z_0 = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{1.42 - 1.5}{0.7/\sqrt{20}} = -0.51$$

$$3. \text{P-value} = P(Z \leq -0.51) = 1 - 0.695 = 0.305$$

## Five – Step Procedure



<eg> pre: #sample 30,  $\mu = 8.16$ ,  $\sigma = 0.17$ ,

CUR:  $\mu = 8.21$

1% level of significance,

$H_0: \text{ave}(\text{cur}) < \text{ave}(\text{pre})$

$$1. \begin{cases} H_0: \mu \geq 8.21 \\ H_1: \mu < 8.21 \end{cases}$$



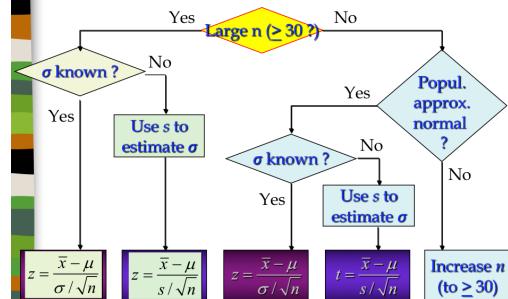
$$2. \alpha = 0.01$$

$$3. \text{TS: } Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

4. Rej.  $H_0$  if  $Z < -2.33$

$$5. Z = \frac{8.16 - 8.21}{0.17/\sqrt{30}} = -2.81 < -2.33 \Rightarrow \text{Rej. } H_0!$$

## Summary of Test Statistics to be Used in a Hypothesis Test about a Population Mean



8

## About 2 Population

## About 2 Population

CI

Sampling Dist of  $\bar{X} - \bar{Y}$

$$\text{let } \bar{X} - \bar{Y} = \bar{Z}$$

$$E\bar{Z} = \mu_1 - \mu_2$$

$$\sigma_{\bar{Z}}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

When  $\sigma_1, \sigma_2$  are Known:

$$\bar{X} - \bar{Y} \pm Z_{\frac{\alpha}{2}} \cdot \sigma_{\bar{X} - \bar{Y}} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Unknown  $\sigma_1, \sigma_2$ , large Sample

$$\bar{X} - \bar{Y} \pm Z_{\frac{\alpha}{2}} \cdot S_{\bar{X} - \bar{Y}} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Unknown  $\sigma = \sigma_1 = \sigma_2$ , Small Sample

$$\bar{X} - \bar{Y} \pm t_{(\frac{n_1+n_2-2}{2}, 0.025)} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \cdot S_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}$$



$$\frac{(\bar{X} - \bar{Y}) - \mu}{\sigma} \sim N(0, 1)$$

$$\Rightarrow \bar{Z} = (\bar{X} - \bar{Y}) \pm \sigma \cdot Z_{\frac{\alpha}{2}}$$

$$= \bar{X} - \bar{Y} \pm Z_{\frac{\alpha}{2}} \cdot \sigma$$

*e.g.*  $n_1 = 120, \bar{x}_1 = 235, s_1 = 15$   
 $n_2 = 80, \bar{x}_2 = 218, s_2 = 20$ , 95% CI

Unknown  $\sigma$ , Large Sample

$$CI = \bar{X} - \bar{Y} \pm Z_{\frac{\alpha}{2}} \cdot S, S = \sqrt{\frac{s_1^2 + s_2^2}{n_1 + n_2}}$$

$$CI = 17 \pm 1.96 \sqrt{\frac{15^2 + 20^2}{120 + 80}}$$

$$= (11.86, 22.14) *$$

*e.g.*  $n_1 = 12, \bar{x}_1 = 29.8, s_1 = 2.56$   
 $n_2 = 8, \bar{x}_2 = 27.3, s_2 = 1.81$ , 95%

Unknown  $\sigma$ , Small Sample

$$CI = \bar{X} - \bar{Y} \pm t_{(n_1+n_2-2, 0.025)} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \cdot S_p$$

$$S_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}$$

$$CI = 29.8 - 27.3 \pm t_{(18, 0.025)} \cdot \sqrt{\frac{1}{12} + \frac{1}{8}} \cdot \sqrt{\frac{11 \cdot 2.56^2 + 7 \cdot 1.81^2}{18}}$$

$$= (0.297, 4.703) *$$

## Hypothesis

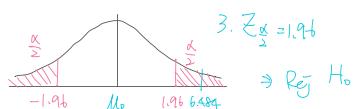
*e.g.* Par Inc.  $n_1 = 120, \bar{x}_1 = 235, s_1 = 15$   
 Rap Inc.  $n_2 = 80, \bar{x}_2 = 218, s_2 = 20$ , 0.05 sign

Is the mean driving distance of Par Inc. golf balls is Greater than the mean driving distance of Rap. Inc. golf balls?

$$1. H_0: \mu_2 - \mu_1 \leq 0$$

$$H_1: \mu_2 > \mu_1$$

$$2. TS = \frac{(\bar{X}_2 - \bar{X}_1 - \mu_2 + \mu_1)}{\sqrt{\frac{s_1^2 + s_2^2}{n_1 + n_2}}} = \frac{(235 - 218) - 0}{\sqrt{\frac{15^2}{120} + \frac{20^2}{80}}} = 6.484$$



$$3. Z_{\frac{\alpha}{2}} = 1.96$$

$$\Rightarrow \text{Rej } H_0$$

## 2 - popu. Problem

### 1 - popu.

Est Mean, Mean<sub>2</sub>

$\mu_1, \mu_2$ : popu. Mean

$\bar{X}_1, \bar{X}_2$ : samp. Mean

Sample Dist:

$$\left\{ \begin{array}{l} \frac{\bar{X}_1 - \bar{X}_2 - \mu_d}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim N(0,1), n \text{ is large} \\ \frac{\bar{X}_1 - \bar{X}_2 - \mu_d}{\sqrt{\frac{s_p^2(\frac{1}{n_1} + \frac{1}{n_2})}} \sim t(n_1+n_2-2), n \text{ is small}} \\ S_p = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \text{ (Weighted Average)} \end{array} \right.$$

$n$  is small:

$$\begin{aligned} \textcircled{1} \quad & n_1, n_2 \text{ are } N\text{-Dist} \\ \textcircled{2} \quad & \sigma_1^2 = \sigma_2^2 = \sigma^2 \\ & S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \end{aligned}$$

$\therefore$  want to plot on the same graph but X-axis only has one dim. So need let  $\sigma_1 = \sigma_2$

$$\frac{\bar{X}_1 - \bar{X}_2 - \mu_d}{\sqrt{S_p^2(\frac{1}{n_1} + \frac{1}{n_2})}} \sim t(n_1+n_2-2)$$

$$\text{Sample Dist: } \begin{cases} \frac{\bar{X}-\mu}{s/\sqrt{n}} \sim Z(0,1) & , n \text{ is large} \\ \frac{\bar{X}-\mu}{s/\sqrt{n}} \sim t(n-1) & , n \text{ is small} \end{cases}$$

$n$  is large

if  
 CLT:  $\begin{cases} \bar{X}_1 \sim N(\mu_1, \frac{s_1^2}{n_1}) \\ \bar{X}_2 \sim N(\mu_2, \frac{s_2^2}{n_2}) \\ \bar{X}_1 \perp \bar{X}_2 \\ \bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}) \\ \frac{\bar{X}_1 - \bar{X}_2 - \mu_d}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim N(0,1) \end{cases}$

(Weighted Average)

<eg>		Sample #1	#2
Inc	Par Inc.		Rap Inc.
# Sample		$n_1=120$	$n_2=80$
Mean		$\bar{X}_1=235$	$\bar{X}_2=218$
SD		$S_1=15$	$S_2=20$
		$\frac{15}{\sqrt{120}} = 1.5$	$\frac{20}{\sqrt{80}} = 2.0$
		$\bar{X}_1 - \bar{X}_2 = 17$	$\pm \sqrt{1.5^2 + 2.0^2} = 5.14$

$\rightarrow$  if  $\theta \neq 0$  & CI

Therefore we can say:  $\Rightarrow \mu_1, \mu_2$  : sig. diff.

I have 95% Conf. to say that My product will better than year in test at

$$S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} = S_p^2$$

$$\bar{X}_1 - \bar{X}_2 \pm t_{0.025} \sqrt{\frac{\bar{X}_1 - \bar{X}_2 - \mu_d}{S_p^2(\frac{1}{n_1} + \frac{1}{n_2})}}$$

Hypo.

$$H_0: \mu_d \geq \mu_0 \quad H_0: \mu_d \leq \mu_0$$

$$H_1: \mu_d < \mu_0 \quad H_1: \mu_d > \mu_0$$

Not useful  $\Rightarrow$  check  $\bar{X}$  in CI or not

$$\boxed{\begin{array}{l} H_0: \mu = 0 \\ H_1: \mu \neq 0 \end{array}}$$

is totally OK!

$$\text{Rej. Regn: } |Z| = \left| \frac{\bar{X}_1 - \bar{X}_2 - \mu_d}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \right| > Z_{\frac{\alpha}{2}}$$

$\Leftrightarrow Z \in CI \Leftrightarrow \text{Accept } H_0$

<eg>		Sample #1	#2
Inc	Par Inc.		Rap Inc.
# Sample		$n_1=120$	$n_2=80$
Mean		$\bar{X}_1=235$	$\bar{X}_2=218$

$$H_0: \mu_d \leq 3$$

Rej Region:

$$H_1: \mu_d > 3$$

$$Z > Z_{0.05}$$

SD	" "	" "
	$S_1 = 15$	$S_2 = 20$

$$H_0: \mu_d \leq 0$$

$$H_1: \mu_d > 0$$

$$\alpha = 0.05 \Rightarrow Z > 1.645 \text{ Rej. Reg.}$$

for  $H_0: \mu_d = 0, H_1: \mu_d \neq 0 \equiv 95\% \text{ CI}$

$$\text{Leg: } n_1 = 22, p = \frac{10}{22}, \bar{x}_1 = 7.125, \bar{x}_2 = 6.45 \\ S_1^2 = 178, S_2^2 = 181$$

$$1. \begin{cases} H_0: \mu_d \leq 0 \\ H_1: \mu_d > 0 \end{cases}$$

$$2. S_p = \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1+n_2-2} = 0.69$$

$$\Rightarrow \text{if } \begin{cases} H_0: \mu_d = 0 \\ H_1: \mu_d < 0 \end{cases} \\ \hookrightarrow p\text{-value} = 0.072$$

$$3. TS = \frac{0.675}{\sqrt{S_p(\frac{1}{n_1} + \frac{1}{n_2})}} \approx 1.9$$

$$4. P\text{-value: } \text{plt}(20) > (.9)$$

$$S_p = \frac{(n_1-1)\sigma_1^2 + (n_2-1)\sigma_2^2}{n_1+n_2-2}$$

$$M_d = \bar{M}_1 - \bar{M}_2$$

$$\frac{\bar{X}_1 - \bar{X}_2 - M_d}{S_p(\frac{1}{n_1} + \frac{1}{n_2})}$$

9



10

## Analysis of Variance

# Analysis of Variance

## One Way ANOVA

### Overview

Goal: Testing for the Equality of  $K$   $\mu$

when  $K=1 \Rightarrow t\text{-test}$   
 when  $K=2 \Rightarrow 2\text{ mean } t\text{-test}$   
 when  $K \geq 3 \Rightarrow \text{One way ANOVA}$

Structure: Predict Variable + Response Variable

Briefly: Check  $H_0: \mu_1 = \mu_2 = \dots = \mu_K$  vs  $H_1: \text{not all } \mu \text{ are equal}$   $\rightarrow$  at least 2  $\mu$  are diff.

Assumptions:

- 1)  $\forall \mu \in N\text{-dist}$

- 2)  $\forall \sigma^2 \text{ are equal}$

- 3)  $\forall \text{obs are indep.}$

### Data

Treatment	Data	Sample Size	Sample Mean	Sample Var
1	$X_{11}, X_{12}, \dots, X_{1n_1}$	$n_1$	$\bar{X}_1$	$S_1^2$
2	$X_{21}, X_{22}, \dots, X_{2n_2}$	$n_2$	$\bar{X}_2$	$S_2^2$
:	:	:	:	:
K	$X_{K1}, X_{K2}, \dots, X_{Kn_K}$	$n_K$	$\bar{X}_K$	$S_K^2$

$\bar{\bar{X}}$ : overall ave.

How & When: When  $\text{Variability}_{\text{Between}} \gg \text{Variability}_{\text{Within}}$   $\Rightarrow$  Reject  $H_0$

### Sum of Square (SS)



$SS_{\text{Total}}$	$\sum_{i=1}^K \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2$	Total SS
$SS_{\text{Treatment}}$ , $SS_{\text{Models}}$	$\sum_{i=1}^K n_i (\bar{X}_i - \bar{X})^2$	SS between groups

Between	$\sum_{i=1}^k \bar{x}_i^2 - \bar{\bar{x}}^2$	$S_{\text{Within}}$
$SS_{\text{Between}}, SS_{\text{Residuals}}$	$\sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$	$SS_{\text{within groups}}$

$\sum (n_i - 1) s_i^2$

## ANOVA Table

Source of Variance	SS	df.	MS	F
Between	$SS_B$	$k-1$	$MS_B$	$\frac{MS_B}{MS_w}$
Within	$SS_w$	$n-k$	$MS_w$	
Total	$SS_T$	$n-1$		

- 1) If  $H_0$  is True,  $\frac{MS_B}{MS_w}$  dist is F-dist ( $k-1, n-k$ )
- 2)  $MS_E$  is Unbiased est. of  $\sigma^2$
- 3) If  $H_0$  is False,  $MS_B$  will overest.  $\sigma^2$

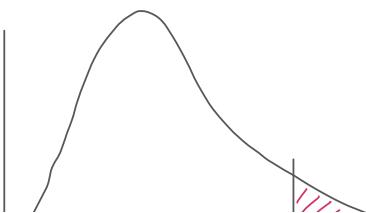
$$H_0 \text{ is } \begin{cases} \text{True : } MS_B \approx MS_w \approx \sigma^2 \\ \text{False : } MS_B \gg MS_w \end{cases}$$

## Hypotheses

$$1) \begin{cases} H_0 : \mu_1 = \mu_2 = \dots = \mu_k \\ H_1 : \text{at least 2 means are not equal.} \end{cases}$$

$$2) TS = F_{(k-1, n-k)} = \frac{MS_B}{MS_w}$$

$$3) \text{Reject Rule : } F > F_\alpha$$



eg.	3 Factories	$\bar{x}$	$s^2$
	Obv1	55	26
	Obv2	68	265
	Obv3	57	24.5

$$1) \begin{cases} H_0 : \mu_1 = \mu_2 = \mu_3 \\ H_1 : \text{not all equal} \end{cases}$$

$$2) T = \frac{MS_B}{MS_w}, \quad MS_B = \frac{SS_B}{k-1}, \quad SS_B = \sum_{i=1}^k n_i (\bar{x}_i - \bar{\bar{x}})^2$$

$$MS_w = \frac{SS_w}{n-k}, \quad SS_w = \sum_{i=1}^k (n_i - 1) s_i^2$$

SS	df	MS	F
490	2	245	
+ +	+ +	+ +	= 9.55
308	12	25.67	
111	14		
798	14		



3)  $F_{0.05}(2, 12) = 3.89 < F \Rightarrow \text{Reject } H_0!$

e.g. Prove  $SS_T = SS_B + SS_W$

e.g.

teaching methods (A, B, C, D). At the level of  $\alpha = 0.01$ , test whether different teaching methods affect score.

Method	Score	$T_i$	$\sqrt{n_i}$	$\chi^2$
A	94, 90, 85, 80	349	4	30516
B	75, 68, 77, 83, 88	391	5	30811
C	70, 73, 76, 78, 80, 68, 65	510	7	37338
D	68, 70, 72, 65, 74, 65	414	6	28634
Total	1664	22		127244
	$\Sigma x$	n		$\Sigma x^2$

$$SS_T = SS_B + SS_W$$

$$SS_T = \sum x - \frac{(\sum x)^2}{n} \quad (1) \dots$$

$$SS_B = \sum \frac{T_i^2}{n_i} - \frac{(\sum x)^2}{n} \quad (2) \quad F = 8.99$$

$$SS_W = SS_T - SS_B \quad (3) \quad F_{0.01}(3, 18) = 5.09 < F \Rightarrow \text{Reject } H_0$$

SS	df	MS	F
890.7	3	296.9	8.99
594.4	18	33	
1485.1	21		

## ANOVA Model

e.g. T-cell Count

$$X_{ij} = \underbrace{\mu}_{\text{Overall popu.}} + \underbrace{\alpha_i}_{\text{Effect}} + \underbrace{\epsilon_{ij}}_{\text{Error Term}}$$

Base level  
Medication  
Unaccounted for Variation

$(\mu = \frac{1}{k} \sum \mu_i)$

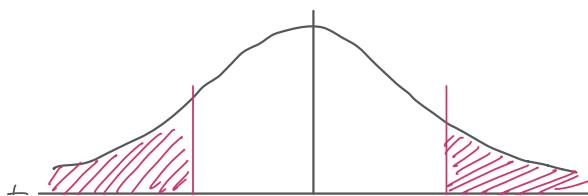
$\left\{ \begin{array}{l} \alpha_1 = \mu_1 - \bar{\mu} \\ \alpha_2 = \mu_2 - \bar{\mu} \\ \alpha_3 = \mu_3 - \bar{\mu} \end{array} \right.$

$\alpha_1 + \alpha_2 + \alpha_3 = 0$

## Fisher's LSD Procedure (check How diff)

1) Hypo:  $H_0: \mu_i = \mu_j$   
 $H_1: \mu_i \neq \mu_j$

2)  $TS = t = \frac{\bar{x}_i - \bar{x}_j}{\sqrt{MS_W (\frac{1}{n_i} + \frac{1}{n_j})}}$



3) Reject Rule:

$$t_{(n-k)} \text{ not in } (-t_{\frac{\alpha}{2}}, t_{\frac{\alpha}{2}})$$

$$2) TS = \bar{x}_i - \bar{x}_j$$

3) Reject Rule

$$|\bar{x}_i - \bar{x}_j| > LSD$$

$$LSD = t_{(\frac{\alpha}{2}, n-k)} \sqrt{MS_w \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}$$

### Example: Reed Manufacturing

#### ■ Analysis of Variance

Observation	Plant 1 Buffalo	Plant 2 Pittsburgh	Plant 3 Detroit
1	48	73	51
2	54	63	63
3	57	66	61
4	54	64	54
5	62	74	56
Sample Mean	55	68	57
Sample Variance	26.0	26.5	24.5

SS	df	MS	F
490	2	245	9.3
308	12	25.67	

$$\begin{cases} H_0: \mu_1 = \mu_2 \\ H_1: \mu_1 \neq \mu_2 \end{cases}$$

TS  $\notin (-6.98, 6.98)$

$$2) TS = |\bar{x}_i - \bar{x}_j| = 13$$

$$3) LSD = t_{(\frac{\alpha}{2}, n-k)} \sqrt{MS_w \left( \frac{1}{n_i} + \frac{1}{n_j} \right)} \\ = 2.18 \cdot \sqrt{25.67 \left( \frac{1}{5} + \frac{1}{5} \right)} = 6.98$$

## Two Way ANOVA

### Overview

Why we need 2-way?

The Predict Variable might be 2 or More

e.g. age, education level  $\rightarrow$  income

### Model

$$x_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$$

Blocking Factor Interested Factor  $\sim N(0, \sigma^2)$

$$\sum \alpha_i = \sum \beta_j = 0$$

### Data

(Block)  
Factor B

	1	2	3	...	b	mean	
(Treatment)	1	$X_{11}$	$X_{12}$	$X_{13}$	$\dots$	$X_{1b}$	$\bar{X}_{1\cdot}$
Factor A	2	$X_{21}$	$\dots$	$\dots$	$\dots$	$\bar{X}_{2\cdot}$	
	3	$X_{31}$	$\dots$	$\dots$	$\dots$	$\vdots$	
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		
	a	$X_{a1}$	$\dots$	$\dots$	$\dots$	$\bar{X}_{a\cdot}$	
mean		$\bar{X}_{\cdot 1}$	$\bar{X}_{\cdot 2}$	$\bar{X}_{\cdot 3}$	$\dots$	$\bar{X}_{\cdot b}$	$\bar{X}$

$$\bar{x}_{ij} - \bar{x} = x_{ij} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x} + (\bar{x}_{i\cdot} - \bar{x}) + (\bar{x}_{\cdot j} - \bar{x})$$

SS

$$\sum_{i=1}^a \sum_{j=1}^b (\bar{x}_{ij} - \bar{x})^2 = \sum_{i=1}^a \sum_{j=1}^b (x_{ij} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x})^2 + b \sum_{i=1}^a (\bar{x}_{i\cdot} - \bar{x})^2 + a \sum_{j=1}^b (\bar{x}_{\cdot j} - \bar{x})^2$$

$$SS_{\text{Total}} = SS_{\text{Error}} + SS_A + SS_B$$

$$df \quad n-1 = (a-1)(b-1) + (a-1) + (b-1)$$

if don't care  
about one factor  
(Block factor)  
then it will be  
Included in  $SS_{\text{Error}}$

### Hypotheses

$$1) \begin{cases} H_{0A} : \text{all } \alpha_i = 0 \\ H_{1A} : \text{not all } \alpha_i = 0 \end{cases}$$

$$2) TS = F_{(df_1, df_2)} = \frac{MS_A}{MSE}$$

3) Reject Rule:

$$F > F_{\alpha}(df_1, df_2)$$

$$1) \begin{cases} H_{0B} : \text{all } \beta_i = 0 \\ H_{1B} : \text{not all } \beta_i = 0 \end{cases}$$

$$2) TS = F_{(df_1, df_2)} = \frac{MS_A}{MSE}$$

3) Reject Rule:

$$F > F_{\alpha}(df_1, df_2)$$

...

## Table

	SS	df	MS	F
Treatment	$\cancel{SS_T}$	$(a-1)$	$\frac{\cancel{SS_T}}{\cancel{df_T}}$	$\frac{MS_T \frac{(a-1)}{B-1}}{MS_E}$
Block	$\cancel{SS_B}$	$(b-1)$	$\frac{\cancel{SS_B}}{\cancel{df_B}}$	$\frac{MS_B \frac{(B-1)}{a-1}}{MS_E}$
Error	$\cancel{SS_E}$	$(a-1)(b-1)$		
Total	Sum [:,1]	Sum [:,2]		
	Sum(SS)	Sum(df)		

eg.

- Eg. There are four routes A, B, C, and D from the suburbs to the downtown area. Is there any difference in driving time? ( $\alpha = 0.05$ )

Route Drivers.	A	B	C	D	$B_r$
John	18	20	20	22	80
Marry	21	22	24	24	91
Anthony	20	23	25	23	91
Claire	25	21	28	25	99
Yvonne	26	24	28	25	103
$T_c$	110	110	125	119	464
$\Sigma x^2$	2466	2430	3169	2839	10904

One way  
Don't care about the Driver effect

ANOVA Table				
Source of variance	Sum of square	df	Mean of square	F
Treatments	32.4	3	10.8	1.618
Error	106.8	16	6.675	
Total	139.2	19		

Two Way

Two-way ANOVA Table				
Source of variance	Sum of square	df	Mean of square	F
Treatments Routes ABD	32.4 (SSTR)	3 (a-1)	10.8 (MSTR)	4.5
Block $SS_{BL}$	78.2 (SSBL)	4 (b-1)	19.6 (MSBL)	4.5
Error $SS_E$	28.6 (SSE)	12 (a-1)(b-1)	2.38 (MSE)	1.83
Total	139.2 (SS <sub>total</sub> )	19 (n-1)		

Where a : number of treatments.  
b : number of blocking.

$SSBL + SSE = \text{OLD } SSE$  (SSE in the previous ANOVA table)

$$\text{Treatment Effect : } F_t = \frac{MS_T}{MS_E} (a-1, (a-1)(b-1)) = 4.5 > 3.49 \Rightarrow \text{Rej } H_{00}$$

$$\text{Block Effect : } F_b = \frac{MS_B}{MS_E} (b-1, (a-1)(b-1)) = 4.5 > 3.259 \Rightarrow \text{Rej } H_{00}$$

## Conclusion

Experiment

- Factor / Treatment Variable : age , education

- Level / Subcategories :  $\left\{ \begin{array}{l} 1-20 \\ 21-40 \\ 41-60 \end{array} \right. , \left\{ \begin{array}{l} \text{Junior} \\ \text{Senior} \\ \text{College} \end{array} \right.$

11

## Chi-Squared Goodness of Fit Test

# Chi-Squared Goodness of Fit Test

## Goodness of Fit Test

### A Multinomial Population

- Categorical Variable : Using Number to token the types
- Binomial - Dist :

Do n times, Succ: X, Fail: n-X

$$X \sim B(n, p) \quad \left\{ \begin{array}{l} E(X) = np \\ \text{Var}(X) = npq, p = \frac{X}{n} \end{array} \right.$$

### Multinomial - Dist

Do n times, type<sub>i</sub>: P<sub>i</sub>, type<sub>j</sub>: P<sub>j</sub> ..., type<sub>k</sub>: P<sub>k</sub>

$$\begin{aligned} f(x_1, x_2, \dots, x_k, p_1, p_2, \dots, p_k) &= \Pr(\bigcap_{i=1}^k X_i = x_i) \\ &= \begin{cases} \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} \cdot p_2^{x_2} \cdots p_k^{x_k}, & \text{when } \sum_{i=1}^k x_i = n \\ 0, & \text{o.w.} \end{cases} \end{aligned}$$

when k=2,  $f(X_1=x_1, X_2=x_2)$   
 $= \frac{n!}{x_1! x_2!} p_1^{x_1} \cdot p_2^{x_2}$   
 $= \binom{n}{x_1} p_1^{x_1} (1-p_1)^{n-x_1}$   
 $X_1 \sim B(n, p_1)$

### Properties

- $E(X_i) = np_i$
- $\text{Var}(X_i) = np_i q_i$
- $\text{Cov}(X_i, X_j) = -np_i p_j$

$$= E(X_i - \bar{X}) - E(X_i)E(\bar{X})$$

$$\begin{aligned} &\text{if } X_i + X_j \sim B(n, p_i + p_j) \\ &\Rightarrow \text{Var}(X_i + X_j) = \text{Var}(X_i) + \text{Var}(X_j) + 2\text{Cov}(X_i, X_j) \\ &\Rightarrow n(p_i + p_j)(1-p_i - p_j) = np_i q_i + np_j q_j + 2\text{Cov}(X_i, X_j) \\ &\Rightarrow \text{Cov}(X_i, X_j) = -np_i p_j \end{aligned}$$

### Hypotheses

- $H_0$ : any
- $H_1$ : not any

Why  $\chi^2$ ?

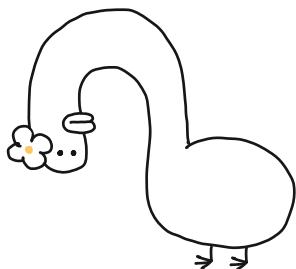
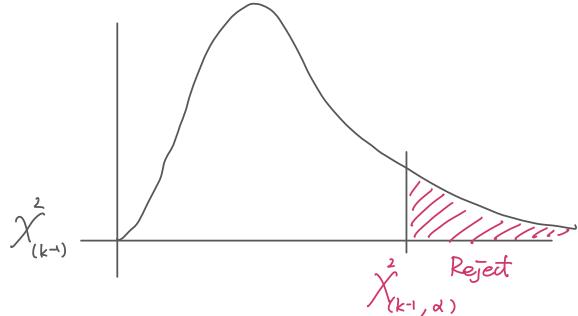
$$\text{Consider } k=2 \Rightarrow \begin{cases} O_1 = X_1 \sim B(n, p_1), E(X_1) = np_1 \\ O_2 = X_2 = n - X_1, E(X_2) = n - np_1 \end{cases}$$

$$\frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} = \frac{(X_1 - np_1)^2}{np_1} + \frac{(X_2 - np_1)^2}{np_1}$$

$$2. TS = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{(k-1)}$$

3. Reject Rule:

$$TS > \chi^2_{\alpha}(k-1)$$



$$E_1 \quad E_2 \quad \frac{n p \mu}{S^2} \approx N(0, 1)^2$$

e.g.

Determine if a six-sided die is fair,

10% sign.level,  $n = 60$

Die Value	1	2	3	4	5	6
Obs. Freq.	9	15	9	8	6	13

$$1. \begin{cases} H_0: P_1 = P_2 = \dots = P_6 = \frac{1}{6} \\ H_1: P_1, \dots, P_6 \text{ are not equal to } \frac{1}{6} \end{cases}$$

$$2. TS = \chi^2_5 \left( \sum_{i=1}^6 \frac{(O_i - E_i)^2}{E_i} \right) = \chi^2_5 \left( \frac{(9-10)^2}{10} + \frac{(15-10)^2}{10} + \dots + \frac{(13-10)^2}{10} \right) = 5.6$$

$$3. TS < \chi^2_{5, 0.1} = 9.24 \Rightarrow \text{do not rej } H_0$$

e.g.

side 2,5, is twice that of other point, 95%

Die Value	1	2	3	4	5	6	Total
Obs. Freq.	13	33	14	7	36	17	120
Exp. Freq.	15	30	15	15	30	15	120

$$1. \begin{cases} H_0: P_2 = P_5 = \frac{2}{8}, P_1 = P_3 = P_4 = P_6 = \frac{1}{8} \\ H_1: \text{not } H_0 \end{cases}$$

$$2. TS = \chi^2 = \frac{1}{8} ((33-30)^2 + (36-30)^2) + \frac{1}{15} ((13-15)^2 + \dots + (17-15)^2) = 6.4$$

$$3. \text{Rj Rule: } TS = 6.4 < 11.07 = \chi^2_{(5, .05)} \Rightarrow \text{Do not reject } H_0$$

## Contingency Table Analysis (Test Independence)

Independence vs. Non-Independence

	D	$\neg D$	Total
M	16	64	80
F	24	96	120
Total	40	160	200

	D	$\neg D$	Total
M	24	56	80
F	6	114	120
Total	30	170	200

$$P(D|M) = 20\%$$

$$P(D|\bar{M}) = 20\%$$

Disease  $\perp$  Sex

$$P(D|M) = 30\%$$

$$P(D|\bar{M}) = 5\%$$

Disease  $\not\perp$  Sex

### Contingency Table

	D	$D'$	Total
M	16 $P_{11}$ $\frac{16}{200}$	64 $P_{21}$ $\frac{64}{200}$	80 $P_{1+}$ $\frac{80}{200}$
F	24 $P_{12}$ $\frac{24}{200}$	96 $P_{22}$ $\frac{96}{200}$	120 $P_{2+}$ $\frac{120}{200}$
Total	40 $P_{+1}$ $\frac{40}{200}$	160 $P_{+2}$ $\frac{160}{200}$	200

$\perp$ : Joint prob. = product of marginal prob.

$$\begin{cases} (M, D) = (M, T) * (T, D) \\ (M, D') = (M, T) * (T, D') \\ (F, D) = (F, T) * (T, D) \\ (F, D') = (F, T) * (T, D') \end{cases}$$

### Test of Independence

Obs. Freq.

	D	$\neg D$	Total
M	$O_{11}$	$O_{12}$	$R_1$
F	$O_{21}$	$O_{22}$	$R_2$
Total	$\sim$	$\sim$	$n$

Exp. Freq. (Indep.)

	D	$\neg D$	Total
M	$E_{11}$	$E_{12}$	$R_1$
F	$E_{21}$	$E_{22}$	$R_2$
Total	$\sim$	$\sim$	$n$

total	$C_1$	$C_2$	$\dots$	$n$
-------	-------	-------	---------	-----

obs. freq.:  $O_{ij}$

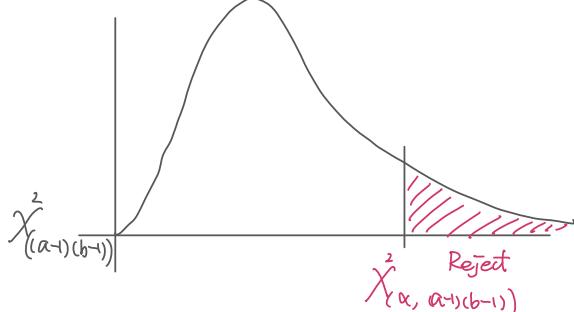
total	$C_1$	$C_2$	$\dots$	$n$
-------	-------	-------	---------	-----

exp. freq.:  $E_{ij} = \frac{1}{n} R_i C_j$

### Hypotheses

1.  $H_0$ : all var. are indep.
1.  $H_1$ : at least 2 var are depen.
2.  $E_{ij} = \frac{1}{n} R_i C_j = \frac{(Total\ row_i) \times (Total\ col_j)}{\text{Sample size}}$
3.  $TS = \sum_{i=1}^a \sum_{j=1}^b \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \sim \chi^2$
4. Reject Rule:

$$TS > \chi^2_{(\alpha, (a-1)(b-1))}$$



eg.

	Obs. Freq.	95%	
	C	NC	Sum(Rn)
M	24	56	80
F	16	104	120
% (Cal)	40	160	200

2. 

	C	NC	% (Rn)
M	40.04 <sup>16</sup>	160.04 <sup>64</sup>	0.4
F	40.06 <sup>24</sup>	160.06 <sup>96</sup>	0.6
% (Cal)	0.2	0.8	1

$$3. TS = \frac{(24-16)^2}{16} + \frac{(56-64)^2}{64} + \frac{(16-24)^2}{24} + \frac{(104-96)^2}{96} = 8.33$$

$$4. TS = 8.33 > \chi^2_{(\alpha, (2-1)(2-1))} = 3.8$$

$\Rightarrow \text{Reject } H_0$

### Does the Data Come From Poisson-Dist?

#### Hypotheses (Poisson)

1.  $H_0$ : data come from Poisson-Dist
1.  $H_1$ : data do not :

2. select RS and Record Obs. freq. ( $f_i$ )
3. Compute Exp. freq. ( $E_i$ ) merge null

eg.

A RS:  $n=100$ , 1-min time intervals recorded  
in the customer arrivals listed below. Is  $X \sim \text{Poi}$ ?

Arrivals	0	1	2	3	4	5	6	7	8	9	10	11	12
freq.	0	1	4	10	14	20	12	12	9	8	6	3	1

$\$ X = \# \text{customers arrival in 1-min.}$

x Unknown

$$4. TS = \sum_{i=1}^k \frac{(f_i - E_i)^2}{E_i} \sim \chi^2_{(k-2)}$$

when  $E_i < 5$

5. Reject Rule:

$$TS > \chi^2_{(\alpha, k-2)}$$

!!! Generally,

$$df = k - p - 1 \quad \begin{cases} K: \# \text{ categories} \\ p: \# \text{ popu. params. est.} \end{cases}$$

$$1. \begin{cases} H_0: X \sim Po(\underline{\mu}) \\ H_1: X \neq Po(\mu) \end{cases}$$

[let  $\mu = \bar{\mu}$ : sample mean]

$$2. \mu = \frac{\text{Total arrivals}}{\text{Total time}} = \frac{\sum \text{Arr}_i \cdot \text{freq}_i}{n} = \frac{600}{100}$$

$$\Rightarrow f(x) = \frac{\mu^x e^{-\mu}}{x!}$$

Arrivals	0	1	2	3	4	5	6	7	8	9	10	11	12
$f(x)$	.0025	.0149	.0446	.0892	.1339	.1620	.1606	.1389	.1041	.0693	.0417	.0227	.0155
$100 \cdot f(x)$	2.5	14.9	44.6	89.2	133.9	162.0	160.6	138.9	104.1	69.3	41.7	22.7	15.5

3. Exp

Arrivals	0	1	2	3	4	5	6	7	8	9	10	11	12
$f(x)$	.0025	.0149	.0446	.0892	.1339	.1620	.1606	.1389	.1041	.0693	.0417	.0227	.0155
$100 \cdot f(x)$	2.5	14.9	44.6	89.2	133.9	162.0	160.6	138.9	104.1	69.3	41.7	22.7	15.5

4. TS

Arrivals	0	1	2	3	4	5	6	7	8	9	10	11	12
$f_i$	5			10	~~~~~	8					10		
$E_i$	6.2			8.92	~~~~~	6.94					7.99		
$f_i - E_i$	-1.2			1.08	~~~~~	-1.06					2.01		

$$K=9!! \quad TS = \frac{(-1.2)^2}{6.2} + \dots + \frac{(-1.06)^2}{7.99} = 3.42 \sim \chi^2_9$$

5. Rej Rule:

$$TS = 3.42 < \chi^2_{(0.05, 9-2)} = 14.07$$

⇒ Do not reject  $H_0$

Hypotheses (Normal Dist)

$$1. \begin{cases} H_0: \text{data come from Normal Dist } (\underline{\mu}, \sigma^2) \\ H_1: \text{data do not } \end{cases}$$

2. select RS and Record Obs. freq. ( $f_i$ )

3. Compute  $E_{\text{exp freq.}}(E_i)$ , merge cell when  $E_i < 5$

$$4. TS = \sum_{i=1}^k \frac{(f_i - E_i)^2}{E_i} \sim \chi^2_{(k-3)}$$

5. Reject Rule:

$$TS > \chi^2_{(\alpha, k-3)}$$

!!! Generally,  $p=2$  since, we have  $\mu, \sigma^2$  to est.

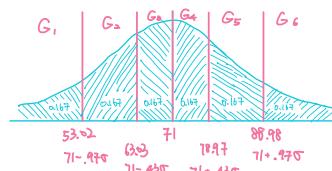
e.g. A RS:  $n=30$ ,  $\bar{x}=71$ ,  $s=18.54$ ,  $N \sim \text{Dist. ?}$

33 43 44 45 52 32 56 58 63 64  
64 65 66 68 70 72 73 73 74 75  
83 84 85 86 91 92 94 98 102 105

$$1. \begin{cases} H_0: X \sim N(\mu, \sigma^2) \\ H_1: X \neq N(\mu, \sigma^2) \end{cases}$$

2. we have  $n=30$ , merge 5 to 1  
therefore, we have 6 intervals →  $K=6$

and each area is  $\frac{1}{6} = 0.167$



group 1: (-∞, 53.02)  
group 2: (53.02, 63.03)

$$df = \underbrace{k - p - 1}_{\text{group}_i} \quad \left\{ \begin{array}{l} K: \# \text{ categories} \\ P: \# \text{ popu. params. est.} \end{array} \right.$$

⋮  
group<sub>i</sub>: (88.98, +∞)

3.  $\curvearrowleft^{K=6}$

group <sub>i</sub>	$f_i$	$E_i$	$f_i - E_i$
<53.02	6	1	1
53.02 ~ 63.03	3	2	-2
63.03 ~ 71	6	5	1
71 ~ 78.97	5	5	0
78.97 ~ 88.98	4	4	-1
>88.98	6	1	1
Total	30	30	

$$4. TS = \sum \frac{(f_i - E_i)^2}{E_i} = 1.6$$

5. Reject Rule:  $TS = 1.6 < \chi^2_{(0.05, 4)} = 7.81$   
 $\Rightarrow$  Do not Reject  $H_0$

12

## Simple Linear Regression

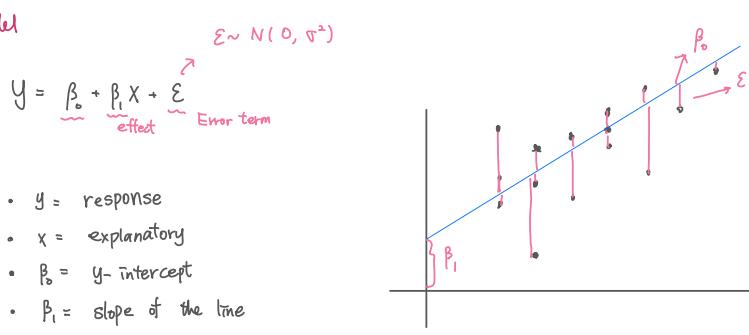
# Simple Linear Regression

## Least Square Procedure

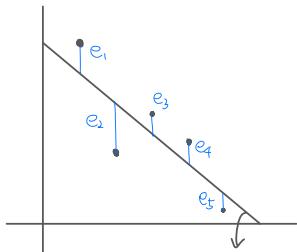
### Target

Examine the relationship between  
explanatory var. ( $x$ ) and response var. ( $y$ )

### Model



### Procedure



1) Assume  $\varepsilon_i \sim N(0, \sigma^2)$

$\varepsilon_i$  are indep.

2) def. SSE (Sum of Squared Errors)

$$\begin{aligned} SSE &= \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 \\ &= \sum (y_i - \beta_0 - \beta_1 x_i)^2 \end{aligned}$$

$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$  3) Find  $\beta_0, \beta_1$  st. SSE has min.

$$\begin{aligned} 4) \quad y &= \hat{\beta}_0 + \hat{\beta}_1 x, \quad \text{where} \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = R_{xy} \cdot \frac{S_y}{S_x} \\ &\Rightarrow y = \bar{y} + \hat{\beta}_1(x - \bar{x}) \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

therefore the CM. is  $(\bar{x}, \bar{y})$

since if  $x = \bar{x}$ , then  $y = \bar{y}$

For minimize SSE

$$\begin{cases} \frac{\partial}{\partial \beta_0} SSE = 0 \\ \frac{\partial}{\partial \beta_1} SSE = 0 \end{cases}$$

we have when  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$  have min SSE

$$\begin{cases} S_{xx} = \sum x_i^2 - \frac{\sum x_i}{n} = \sum (x_i - \bar{x})^2 \\ S_{xy} = \sum x_i y_i - \frac{\sum x_i \sum y_i}{n} = \sum (x_i - \bar{x})(y_i - \bar{y}) \\ S_{yy} = \sum y_i^2 - \frac{\sum y_i}{n} = \sum (y_i - \bar{y})^2 \end{cases}$$

pf>

$$\min \sum_i (y_i - \beta_0 - \beta_1 x_i)^2 \triangleq f(\beta_0, \beta_1)$$

$$\begin{cases} \frac{\partial}{\partial \beta_0} f = -2 \sum_i (y_i - \beta_0 - \beta_1 x_i) = 0 \\ \frac{\partial}{\partial \beta_1} f = -2 \sum_i x_i (y_i - \beta_0 - \beta_1 x_i) = 0 \end{cases}$$

$$\Rightarrow \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

Example: The data  $y$  has been observed for various values of  $x$ , as follows:

y	240	181	193	155	172	110	113	75	94
x	1.6	9.4	15.5	20.0	22.0	35.5	43.0	40.5	33.0

Fit the simple linear regression model using least squares.

e.g.

$$\hat{\beta}_0 = S_{xx} \quad r^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}}$$

$$1. \begin{cases} P_i = \bar{S}_{xy} \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases} \quad \begin{cases} S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n} \\ S_{xy} = \sum x_i y_i - \frac{\sum x_i \sum y_i}{n} \end{cases}$$

$$\begin{aligned} S_{xx} &= \sum (x_i - \bar{x})^2 \\ &= \sum x_i^2 - n\bar{x}^2 \end{aligned}$$

$$= \frac{1}{S_{xx}} \begin{bmatrix} \frac{1}{n} \sum x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

n	$\Sigma x_i$	$\Sigma x_i^2$	$\Sigma y_i$	$\Sigma y_i^2$	$\Sigma x_i y_i$
9	220.5	703.7	1330	2205.49	24864.4

$$3. \begin{cases} S_{xx} = 703.7 - \frac{1}{9}(220.5)^2 = 1651.45 \\ S_{xy} = 24864.4 - \frac{1}{9}(220.5 \cdot 1330) = -5794.1 \end{cases}$$

$$4. \begin{cases} \hat{\beta}_1 = \frac{1651.45}{-5794.1} = -3.5085 \\ \hat{\beta}_0 = \frac{1}{9}(1330) - (-3.5085) \times \frac{1}{9}(220.5) = 234.1 \end{cases}$$

$$5. \Rightarrow y = 234.1 - 3.5085 x$$

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum x_i y_i + \sum x_i^2 \sum x_i y_i}{\sum x_i^2 - n\bar{x}^2} \\ &= \frac{\sum x_i y_i}{S_{xx}} = \frac{S_{xy}}{S_{xx}} \end{aligned}$$

$$\Rightarrow \hat{\beta}_0 = \frac{1}{S_{xx}} \left( \sum x_i^2 \bar{y} - (\sum x_i y_i) \bar{x} \right)$$

$$= \bar{y} - \hat{\beta}_1 \bar{x}$$

stick-up.

Quantifying the Goodness of the fit

residuals

$$\text{Est. error} = \hat{e}_i = y_i - \hat{y}_i$$

usage residuals to Est. error

therefore the  $e_i \approx \frac{1}{n-1} \sum \hat{e}_i^2$  residuals

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n-2} \sum \hat{e}_i^2 = \frac{1}{n-2} \sum (y_i - \hat{y}_i)^2 \\ &= \frac{S_{xx} - 2S_{xy}}{n-2} \left( \frac{\text{Residual Sum of Squares}}{\text{df}} \right) \end{aligned}$$

$\beta_0, \beta_1$  have been est.

$$\hat{\beta}_1 = r_{xy} \cdot \frac{s_y}{s_x}$$

$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$

$$S_{xx} = (n-1)s_x^2 \quad \downarrow \quad r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2}}$$

$$\begin{cases} s_x^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 \\ s_y^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2 \end{cases}$$

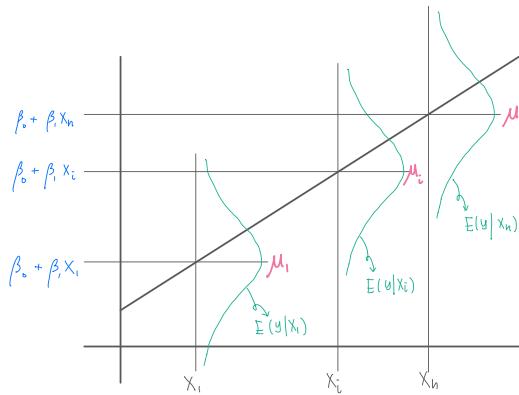
$$\begin{cases} s_x^2 = S_{xx} \\ s_y^2 = S_{yy} \end{cases}$$

Sum of Error of Square (SSE)

SSE

$$\begin{aligned}
 SSE &= \sum_i (y_i - \hat{y}_i)^2 \\
 &= S_{yy} - \hat{\beta}_1 \cdot S_{xy} \\
 &= S_{yy} - \hat{\beta}_1^2 \cdot S_{xx} \\
 &\quad \xrightarrow{\text{pf}} \\
 &= \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\
 &= \sum (y_i - \bar{y} - \hat{\beta}_1 (x_i - \bar{x}))^2 \\
 &= \sum (y_i - \bar{y})^2 - 2 \hat{\beta}_1 \cdot \sum (y_i - \bar{y})(x_i - \bar{x}) \\
 &\quad + \hat{\beta}_1^2 \cdot \sum (x_i - \bar{x})^2 \\
 &= S_{yy} - \hat{\beta}_1 \cdot S_{xy} \\
 \hat{\sigma}^2 &= S_{\epsilon}^2 = \frac{SSE}{n-2}
 \end{aligned}$$

## The Normality of $\epsilon$



## Hypotheses ( $\beta_1$ : slope)

$$\begin{aligned}
 1) \quad & \left\{ H_0: \beta_1 = \beta_1' \text{ or } H_0: \beta_1 = \beta_1'' \right. \\
 & \left. H_1: \beta_1 \neq \beta_1' \quad H_1: \beta_1 > \beta_1' \right. \dots
 \end{aligned}$$

2) check:  $r^2$

3) est. Residuals:  $\hat{e}_i = y_i - \hat{y}_i$

eg.

Global Temperature (cont.)						
SUMMARY OUTPUT						
Regression Statistics						
Multiple R	0.85739864					
R Square	0.735717675					
Adjusted R Square	0.729425238					
Standard Error	0.131613902					
Observations	44					

ANOVA						
Results of hypothesis test						
Regression	$H_0: \beta_1 = 0$	<b>Reject <math>H_0</math></b>				
Residual	$H_A: \beta_1 \neq 0$		42	0.727533203	0.017322219	<b>Confidence interval for <math>\beta_1</math></b>
Total			43	2.752863636		

Coefficients	Standard Error	T Stat	P-value	Lower 95%	Upper 95%
Intercept	10.6441999	0.341822134	31.1395819	1.19922E-30	9.954374809
CO2	0.010324616	0.000954834	10.81299943	1.03354E-13	0.008397684

## Test the Slope ( $X, Y$ relationship)

If Slope ( $\hat{\beta}_1$ )  $\approx 0$ , then  $X, Y$  have no "Linear Relationship."

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \begin{cases} E(Y_i) = \beta_0 + \beta_1 x_i \\ \text{Var}(Y_i) = \text{Var}(\epsilon_i) = \sigma^2 \end{cases}$$

$$\bullet S_{xy} = \sum_i (x_i - \bar{x})(y_i - \bar{y}) = \sum_i (x_i - \bar{x}) y_i$$

$$\bullet \text{Var } S_{xy} = \sum_i (x_i - \bar{x})^2 \cdot \text{Var}(y_i)$$

1.	$H_0: \beta_1 = 0$	$H_1: \beta_1 \neq 0$	$\Rightarrow t^* = S_{xx} \cdot \hat{\beta}_1$
2.	$TS \sim t - D_{ist} = \frac{\hat{\beta}_1 - \beta_1}{S_{\hat{\beta}_1}}$		$\bullet \text{Var } \frac{S_{xy}}{S_{xx}} = \frac{\text{Var } S_{xy}}{S_{xx}^2} = \frac{\sigma^2}{S_{xx}}$
	$S_{\hat{\beta}_1} = \frac{S_\epsilon}{\sqrt{n-1} \cdot S_x}$		$\bullet S_{\hat{\beta}_1} = \sqrt{\text{Var } \hat{\beta}_1} = \sqrt{\frac{S_{xy}}{S_{xx}}} = \frac{\sigma}{\sqrt{S_{xx}}}$
3.	Reject Rule:		$\bullet \hat{\sigma}^2 = \frac{SSE}{n-2}, \frac{S_{xx}}{n-1} = S_{\hat{\beta}_1}^2 = S_x^2$

\*  $CI = \hat{\beta}_1 \pm t_{(n-2, \alpha/2)} S_{\hat{\beta}_1}$

## Coefficient of Determination ( $r^2 = \frac{SS_R}{SS_T}$ )

If  $r^2 \uparrow$ , then Model fixed well!

• Case I) $\begin{cases} \text{Ignored } X \\ \text{use } \bar{y} \text{ to predict } y \end{cases}$	$SS_{\text{Total}} = \sum (\text{obs} - \text{pred})^2$	$SS_T = SS_E + \underbrace{SS_{\text{Regression}}}_{\hat{\beta}_1 S_{xy}} \approx \hat{\beta}_1 S_{xy}$
	$\frac{SS_{\text{Total}}}{S_{yy}} = \sum (y_i - \bar{y})^2$	$\begin{aligned} SS_T &= \sum (y_i - \bar{y})^2 \\ &= \sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\ &= \underbrace{\sum (y_i - \hat{y}_i)^2}_{SS_E} + \underbrace{\sum (y_i - \bar{y})^2}_{SS_R} + 2 \sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \end{aligned}$
• Case II) $\begin{cases} \text{use } X \\ \text{use } \hat{y} = \beta_0 + \beta_1 X \end{cases}$	$SS_{\text{Error}} = \sum (\text{obs} - \text{pred})^2$	$r^2 = \frac{SS_R}{SS_T}, \text{ if } r^2 \approx 1, \text{ then better}$
	$\frac{SS_{\text{Error}}}{S_{yy} - \hat{\beta}_1 S_{xy}} = \sum (y_i - \hat{y}_i)^2$	$\begin{aligned} \frac{SS_R}{SS_T} &= 1 - \frac{SS_E}{SS_T} = 1 - \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{S_{yy}} \\ &= \frac{\hat{\beta}_1 S_{xy}}{S_{yy}} = \frac{S_{xy} S_{xy}}{S_{yy} S_{yy}} \\ &= \frac{S_{xy}^2}{S_{xx} \cdot S_{yy}} = (r_{xy})^2 \end{aligned}$

eg.

	SS	df	MS	F	P-value
Reg.	1065.8	1	1065.8	10.87	0.01
Err.	784.2	8	98		$F_{dist}(1, 8) = 10.87$
Total	1850	9			

$$\begin{aligned} t^* &= \left( \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}} \right)^2 \\ \Rightarrow F &= \frac{MS_R}{MS_E} = \frac{\frac{SS_R}{df_R}}{\frac{SS_E}{df_E}} = \frac{\frac{SS_R}{1}}{\frac{SS_E}{8}} = \frac{SS_R}{8} \\ \therefore F &= \frac{1}{8} \cdot \frac{S_{xy}^2}{S_{xx}} \\ \Rightarrow t &= \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}}, t^* = \frac{\hat{\beta}_1^2}{S_{\hat{\beta}_1}^2} = \frac{\left( \frac{S_{xy}}{S_{xx}} \right)^2}{\frac{1}{8} \cdot \frac{S_{xy}^2}{S_{xx}}} \\ \therefore t^* &= \frac{1}{8} \cdot \frac{S_{xy}^2}{S_{xx}} = F \end{aligned}$$

## Using the Regression Equation

Point Estimation → Interval Estimation

when we have Reg Eq ( $\hat{Y} = \beta_0 + \beta_1 X$ )

Just input  $X$ , we could give  $\hat{Y}$  (point)  
But "How Close" to the real point?

Confidence Interval (for "average" of all 95% CI)

Est. the "Average Y" for given  $X_v$

$$\hat{Y} \pm t_{\frac{\alpha}{2}} \cdot \sqrt{SE^2(\hat{\beta}_1) \cdot (X_v - \bar{X})^2 + \frac{S_e^2}{n}}$$

df: n-2

Predict Interval (for "this"  $X_v$ 's prediction)

"Predict Y" for given  $X_v$

$$\hat{Y} \pm t_{\frac{\alpha}{2}} \cdot \sqrt{SE^2(\hat{\beta}_1) \cdot (X_v - \bar{X})^2 + \frac{S_e^2}{n} + \frac{S_e^2}{\hat{\sigma}^2}}$$

C.I.

$$y \mid X_v \sim N(\beta_0 + \beta_1 X_v, \hat{\sigma}^2)$$

$$\Rightarrow E(y \mid X_v) = \hat{\beta}_0 + \hat{\beta}_1 X_v$$

$$\Rightarrow \text{Var}(\hat{E}(y \mid X_v))$$

$$= V_{X_v}(\hat{Y}) + (X_v - \bar{X})^2 \cdot \text{Var}(\hat{\beta}_1) + 2 C_{\alpha}(\hat{Y}, \hat{\beta}_1(X_v - \bar{X}))$$

$$= \frac{\hat{\sigma}^2}{n} + (X_v - \bar{X})^2 \cdot \frac{\hat{\sigma}^2}{S_{xx}}$$

P.I.

$$Y^* = \beta_0 + \beta_1 X_v + \varepsilon$$

$\Rightarrow$  Predict  $Y^*$ : C.I.  
though  $\varepsilon = 0$ , but we still should consider the effect to Variation of Error ( $\varepsilon$ )

$$\text{Var} Y^* = \text{Var}(\hat{E}(y \mid X_v)) + \text{Var}(\varepsilon)$$

$$= \frac{\hat{\sigma}^2}{n} + (X_v - \bar{X})^2 \cdot \frac{\hat{\sigma}^2}{S_{xx}} + S_{\varepsilon}^2$$

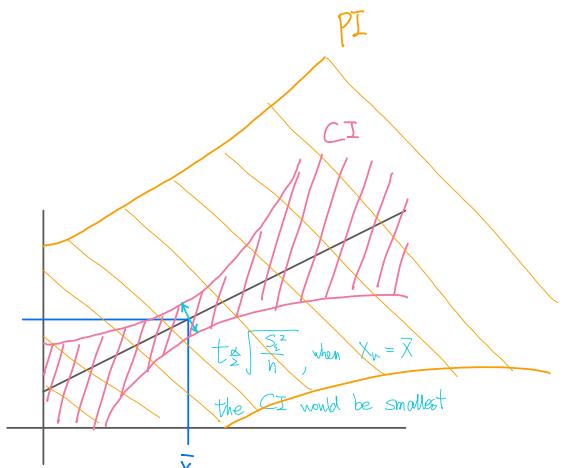
$X_v$ 's effect to the Interval

If  $|X_v - \bar{X}| \uparrow$  then CI ↑

- $CI = \hat{Y} \pm t_{\frac{\alpha}{2}} \cdot \sqrt{SE^2(\hat{\beta}_1) \cdot (X_v - \bar{X})^2 + \frac{S_e^2}{n}}$

If  $X_v = \bar{X}$ , then  $CI = \hat{Y} \pm t_{\frac{\alpha}{2}} \sqrt{0 + \frac{S_e^2}{n}}$

- $PI = \hat{Y} \pm t_{\frac{\alpha}{2}} \sqrt{SE^2(\hat{\beta}_1) \cdot (X_v - \bar{X})^2 + \frac{S_e^2}{n} + S_{\varepsilon}^2}$



## Regression Diagnostics

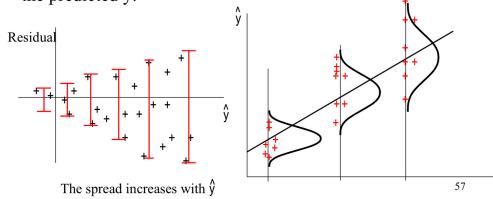
$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

*y is correct?*

- $\varepsilon_i$  is N-Distr? ( $\varepsilon_i \sim N()$ ?)
- $\forall x_i, \varepsilon_i \in \text{const}$ ? ( $x_i \perp \varepsilon_i$ ?)
- $\forall \varepsilon_i$  are indep.? ( $\perp_{\text{iid}} \varepsilon_i$ ?)

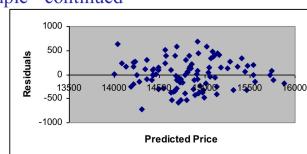
## Residual Analysis

- When the requirement that  $SD(\varepsilon)$  is the same for all  $x$  is violated, we have a condition called heteroscedasticity.
- Diagnose heteroscedasticity by plotting the residual against the predicted  $\hat{y}$ .



### Homoscedasticity

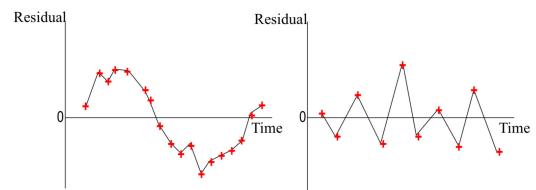
- When the requirement that  $SD(\varepsilon)$  is the same for all  $x$  is not violated we have a condition of homoscedasticity.
- Example - continued



## Non-Independence of Error Variables

- A **time series** is constituted if data were collected over time.
- Examining the residuals over time, no pattern should be observed if the errors are independent.
- When a pattern is detected, the errors are said to be autocorrelated.
- Autocorrelation can be detected by graphing the residuals against time.

Patterns in the appearance of the residuals over time indicates that autocorrelation exists.



Note the runs of positive residuals, replaced by runs of negative residuals.

Note the oscillating behavior of the residuals around zero.