

## c1. Descriptive Statistics .

1.1 Graphical Representation .

1.2 Statistical Descriptors

1.2.1 Measure of Central Tendency .

1.2.2 Measure of Dispersion

1.2.3 Measure of Asymmetry

1.2.4 Measure Relations between Datasets

## c2. Probability

2.1 Probability

2.2 Conditional Probability

## c3. Random Variables

3.1 Random Variables

3.2 Function of Random Variables ( Conversion of PMF and PDF )

3.3 Moments

3.4 Expectation Algebra

## c4. Probability Distribution

4.1 Discrete Distribution

4.2 Continuous Distribution

4.3 Derivation of  $E(X)$  and  $\text{Var}(X)$  for all Distribution

## c5. Fitting Probability Distribution

5.1 Method of Moments

5.2 Maximum Likelihood Method

5.3 Goodness of Fits

## c6. Multiple Variables

6.1 Covariance and Correlation .

6.1.1 Covariance

6.1.2 Properties of Covariance

6.1.3 Correlation Coefficient .

6.1.4 Properties of Correlation Coefficient

6.2 Joint Distribution

6.3 Marginal Distribution

6.4 Conditional Distribution

6.5 Important Identity

6.6 Sum of Two Random Variables

6.6.1 Method of Convolution

6.6.2 Moments

6.6.3 Distributions

## 6.7 Bivariate Normal Distribution

### c7. Confidence Interval

7.0 Pre-requisites

7.1 Introduction to  $E(\bar{X})$  and  $\text{Var}(\bar{X})$

7.2 Limit Theorems

7.4 Sampling Distribution

7.4.1 Sampling Distribution of Sample Mean

7.4.2 Sampling Distribution of Sample Variance

### c8 Statistical Testing

8.1 Introduction

8.2 Asymmetric Case (One-sample Analysis) and Symmetric Case (Two-sample Analysis)

8.3 Test Error

8.4 Goodness of Fit Test

8.4.1 Chi-squared Test ( $\chi^2$ )

8.4.2 Kolmogorov-Smirnov (KS) Test

### c9 Linear Regression

9.1 Introduction

9.2 Properties of Parameter Estimates

9.3 Confidence Interval (Population Parameters :  $a$  and  $b$ )

9.4 Confidence Interval ( $E(Y|X=x_0)$  and  $Y_{x=x_0}$ )

9.5 Statistical Testing (Hypothesis Testing)

9.5.1 Slope Parameter Test

9.5.2 Correlation Test

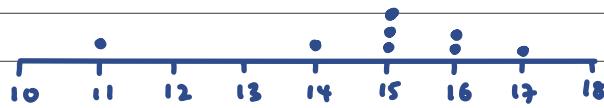
9.5.3 Regression Test (F-test)

## c1. Descriptive Statistics

### 1.1 Graphical Representation.

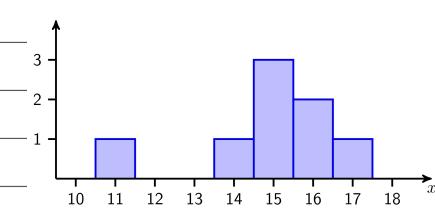
#### 1. Dot plot.

$$\{11, 14, 15, 16, 17, 15, 16, 15\}$$

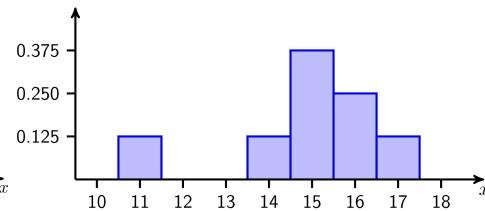


#### 2. Histogram.

$$\{11, 14, 15, 16, 17, 15, 16, 15\}$$

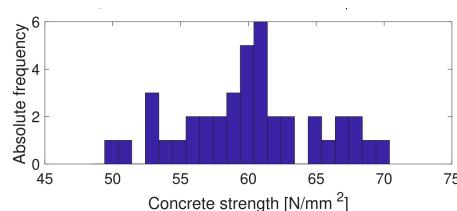


(absolute freq.)

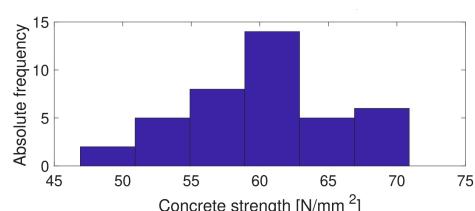


(relative freq.)

} different in frequency type



(bin size = 1)



(bin size = 4)

- bin size does not have to be identical in a single plot!

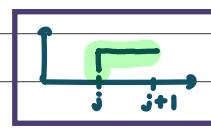
} different in bin size .

#### 3. CDF for discrete (more on CDF in c3)

(cumulative relative frequency)

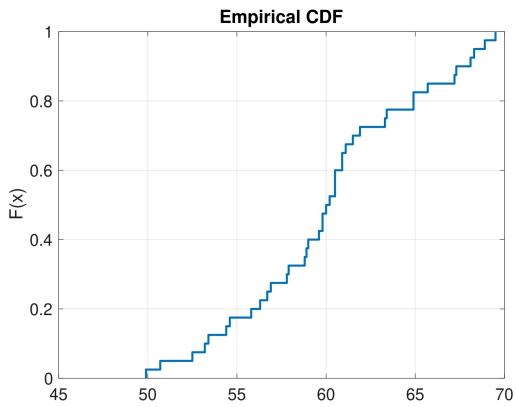
formula is given by:

$$F(x) = \begin{cases} 0 & \text{for } x < x_{(1)} \\ j/n & \text{for } x_{(j)} \leq x < x_{(j+1)} \text{ and } 1 \leq j \leq n-1 \\ 1 & \text{for } x \geq x_{(n)} \end{cases}$$



all the green parts takes all the frequency before  $x_{(j+1)}$

it's like bump up the cumulative frequency only once you reach exactly the  $x$ .



**eg1:**

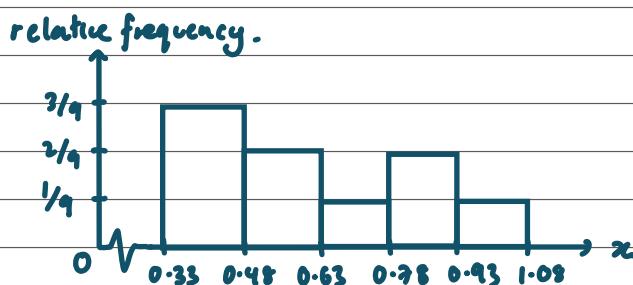
The current speed in the River Thames has been measured every 3 hours over the course of a day in m/s as given below:

$$\{0.54; 0.82; 0.71; 0.92; 0.41; 1.07; 0.41; 0.56; 0.87\}$$

- (a) Draw a histogram of the dataset with 0.15 m/s bin widths, starting at 0.33 m/s.
- (b) Plot the empirical cumulative distribution frequency of the current speed.

**(a)**

Histogram



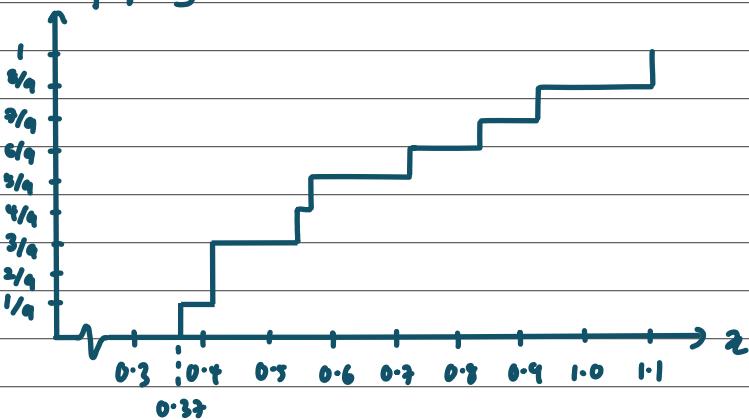
| Bin         | f |
|-------------|---|
| 0.33 - 0.48 | 3 |
| 0.48 - 0.63 | 2 |
| 0.63 - 0.78 | 1 |
| 0.78 - 0.93 | 2 |
| 0.93 - 1.08 | 1 |

$$\sum f = 9, \frac{1}{9} = 0.11$$

**(b)**

CDF

Relative frequency,  $F(x)$



| j | $x_j$ | $F(x_j)$ |
|---|-------|----------|
| 1 | 0.33  | 1/9      |
| 2 | 0.41  | 2/9      |
| 3 | 0.41  | 3/9      |
| 4 | 0.54  | 4/9      |
| 5 | 0.64  | 5/9      |
| 6 | 0.71  | 6/9      |
| 7 | 0.82  | 7/9      |
| 8 | 0.91  | 8/9      |
| 9 | 1.07  | 1        |

$$n=9$$

## ★ How To KNOW TO USE SAMPLE OR POPULATION?

Real data → use sample.

e.g. "given these numbers...", "from this survey..."

### 1.2 Statistical Descriptors.

Theoretical distributions → population.

e.g. "X follows normal distribution..."

#### 1.2.1 Measure of Central Tendency.

∴ Basically if experiment is conducted with real data, use sample.

Describe around which central value the data cluster.

Describe with: mean, mode, median.

##### 1. Mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{\Sigma x}{n}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n f_i x_i = \frac{\Sigma f x}{\Sigma f}$$

##### 2. Median.

$$\text{median} = x_{\left(\frac{n+1}{2}\right)}$$

(if n → odd)

$$\text{median} = \frac{x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)}}{2}$$

(if n → even)

##### 3. Mode.

mode = x with the highest frequency.

(can have more than one mode → multi-modal)

##### 4. Geometric Mean.

$$\hat{x} = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

(geometric mean)

## 1.2.2 Measure of Dispersion.

Describe how spread out is the data.

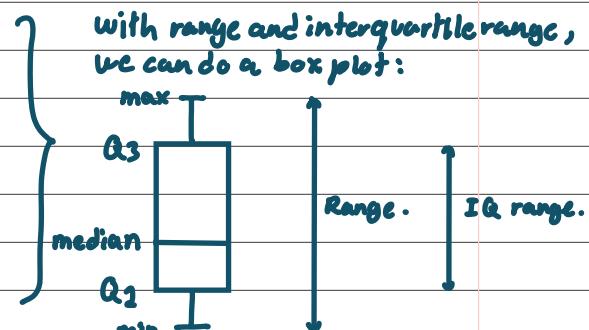
Described with: range, interquartile range, variance, standard deviation, mean absolute deviation, coeff. of variance.

### 1. Range.

$$\text{range} = \max(x) - \min(x) \\ = x(n) - x(1)$$

### 2. Interquartile Range.

$$IQR = q_{75}(x) - q_{25}(x)$$



how to find quantile?

#### Method 1:

for  $q_p(x)$ , find  $j$  such that:

$$\frac{j-0.5}{n} < \frac{p}{100} \leq \frac{j+0.5}{n}$$

and take mean values of  $x_j$  and  $x_{j+1}$ .

e.g.  $x = \{1, 4, 6, 6, 7\}$ , find  $q_{25}(x)$

test  $j=1$ :

$$\frac{1-0.5}{5} = 0.1 ; \frac{1+0.5}{5} = 0.3 \rightarrow 0.1 \leq 0.25 \leq 0.3 \rightarrow j=1 \text{ works!}$$

$$q_{25}(x) = \frac{x(1) + x(1+1)}{2} = \frac{1+4}{2} = 2.5$$

#### Method 2:

find  $k = p/100(n+1)$

if  $k$  lies on an integer, then  $x(k) = q_p(x)$

if  $k$  lies between two numbers (e.g. 2.25), then  $\frac{x(s) + x(s)}{2} = q_p(x)$

3. Variance.  $\sigma^2$  population variance,  $s^2$  sample variance,  $\sigma^2$  population variance is biased low (underestimates true variance)

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 \left( \frac{\sum f x^2}{\sum f} - \left( \frac{\sum f x}{\sum f} \right)^2 \right) \text{ with frequencies version.}$$

$\underbrace{\quad \quad \quad}_{\text{easier to use.}}$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \left. \begin{array}{l} \text{no simplified version formula} \\ \longrightarrow \frac{1}{n-1} \sum f (x - \bar{x})^2 \quad (\text{with frequencies version}) \end{array} \right.$$

4. Standard Deviation.

$$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \quad \left( \sqrt{\frac{1}{n} \sum f (x - \bar{x})^2} \right) \text{ with frequencies}$$

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad \left( \sqrt{\frac{1}{n-1} \sum f (x - \bar{x})^2} \right) \text{ with frequencies}$$

5. Coefficient of Variation.

used to compare the degree of spread across multiple datasets.

$$V = \frac{s}{\bar{x}}$$

6. Mean Absolute Deviation.

Better version of standard deviation.

because less influenced by extreme values.

$$d = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}| = \frac{1}{n} \sum |x - \bar{x}|$$

### 1.2.3 Measure of Asymmetry.

sample skewness can use both of these

1. Biased skewness.

$$g_1(x) = \frac{1}{n} E(x - \bar{x})^3$$

2. Unbiased skewness.

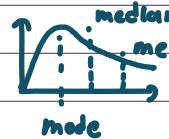
$$\hat{g}_1(x) = \frac{\sqrt{n(n-1)}}{n-2} g_1(x)$$

$$G^3 = \left( \frac{E(x - \bar{x})^3}{n} \right)^{2/3}$$

population skewness.  
currently use this.

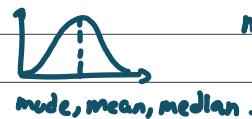
We will use more biased skewness than unbiased.

i.  $g_1(x) > 0$  (positive skewness)



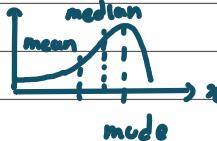
mode < median < mean.

ii.  $g_1(x) \approx 0$  (symmetric)



mode = median = mean

iii.  $g_1(x) < 0$  (negative skewness)



mode > median > mean.

### 1.2.4 Measure Relation between Datasets. (to check whether there is any LINEAR correlation between $x_i$ and $y_i$ )

#### 1. SAMPLE covariance.

(problem about this

is  $\text{cov}_{xy}$  has units.

if  $x$  and  $y$  has unit [m]

then  $\text{cov}_{xy}$  has unit [ $m^2$ ])

$$\text{cov}_{x,y} = \frac{1}{n-1} E(x - \bar{x})(y - \bar{y}) \xrightarrow{\text{so to non-dimensionalise...}}$$

#### 2. SAMPLE correlation coefficient.

$$C_{x,y} = \frac{\text{cov}_{x,y}}{S_x S_y}, \text{ bounded by: } -1 \leq C_{x,y} \leq 1$$

\* if population:

#### 1. Population covariance

$$\text{Cov}(X, Y) = \frac{1}{n} E(x - \mu_x)(y - \mu_y) = E(XY) - E(X)E(Y)$$

i.  $C_{x,y} = -1$ :



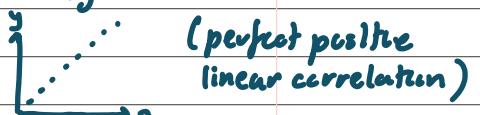
(perfect negative linear correlation)

ii.  $C_{x,y} = 0$



(completely no correlation)

iii.  $C_{x,y} = 1$



(perfect positive linear correlation)

#### 2. Population correlation coefficient.

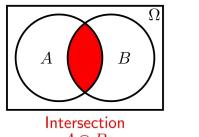
$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{S_x S_y}$$

## C2. Probability

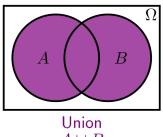
### 2.1 Probability.

#### Venn Diagram

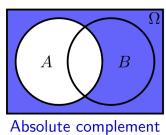
Consider two sets  $A$  and  $B$  in the universal set  $\Omega$



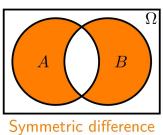
Intersection  
 $A \cap B$



Union  
 $A \cup B$



Absolute complement  
 $\bar{A}$



Symmetric difference  
 $A \Delta B$

#### Types of event.

same row  
(do something once)



1. Mutually Exclusive: events cannot happen at the same time.  $P(A \cap B) = 0$
2. Mutually Inclusive: events that can happen at the same time.  $P(A \cap B) \neq 0$

- {
3. Independent: event will not be affected by any events.  $P(A|B) = P(A)$ ;  $P(B|A) = P(B)$
  4. Dependent: event that will be affected by other events.

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

two rows  
(do two things)

#### 3 Theorems of Probability.

$$1. P(\bar{A}) = 1 - P(A)$$

$$2. \text{ if } A_i \cap A_j = \emptyset \ \forall i, j : P(A_1 \cup A_2 \cup \dots \cup A_k) = \sum_{i=1}^k P(A_i) \rightarrow P(A \cup B) = P(A) + P(B)$$

$$3. \text{ if } A \cap B \neq \emptyset : P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

" $P_k$  is permutation (keyword: in order, sequence is important)

" $C_k$  is combination (keyword: choose, without order, sequence is not important)

$$nP_k = \frac{n!}{(n-k)!} ; {}^nC_k = \frac{n!}{(n-k)!k!}$$



#### Addition Rule

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

if mutually exclusive:

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



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

if independent:

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

Bayes' Law

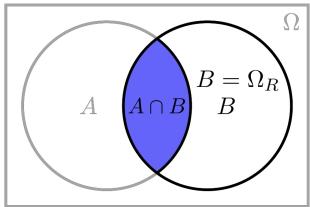
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

## 2.2 Conditional Probability.



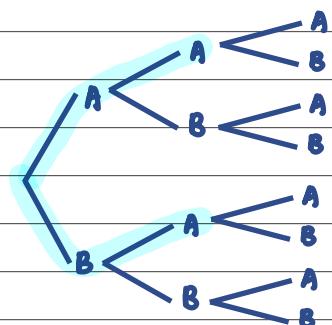
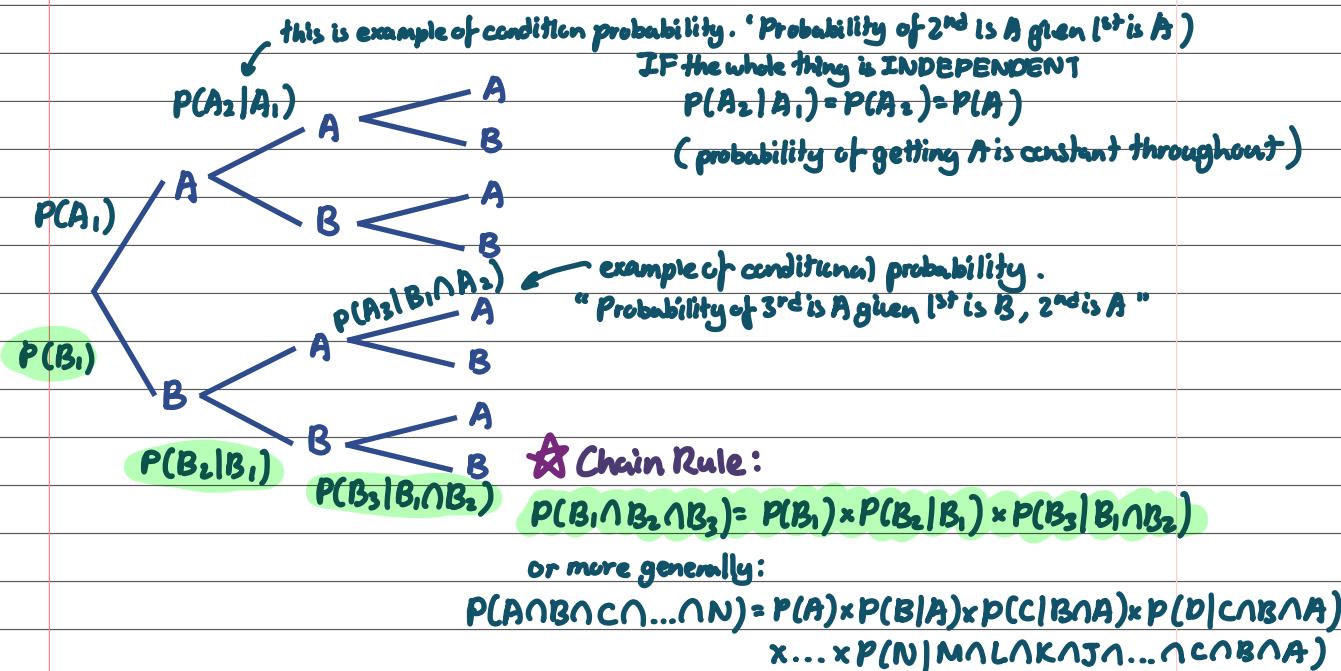
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

"Probability of A given B (already happened)"



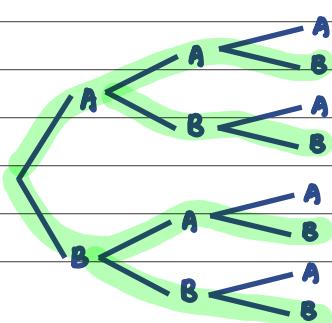
sample space from  $\Omega$  reduced to  $\Omega_R$ .

**TREE DIAGRAM** (I will explain all the rule of Conditional Probability with Tree Diagram)



★ Law of Total Probability :  $P(B) = \sum_{i=1}^k P(B|A_i) P(A_i)$

$$\text{eg1. } P(A_2) = P(A_1 \cap A_2) + P(B_1 \cap A_2) \\ = P(A_2|A_1)P(A_1) + P(A_2|B_1)P(B_1)$$



$$\text{eg2. } P(B_3) = P(A_1 \cap A_2 \cap B_3) + P(B_1 \cap B_2 \cap B_3) \\ + P(B_3|A_1 \cap A_2)P(A_2|A_1)P(A_1) \\ + P(B_3|A_1 \cap B_2)P(B_2|A_1)P(A_1) \\ + \dots$$

chain rule.

## c3. Random Variables.

### 3.1 Random Variables

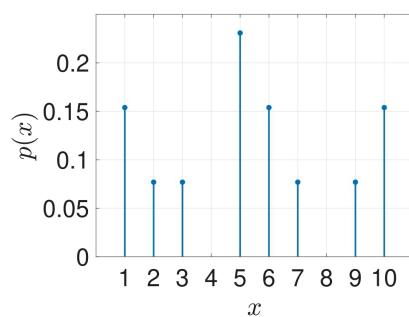
#### Probability Mass Function (PMF)

- tell us the probability of each DISCRETE random variable.

$$\cdot p(x) = P(X=k)$$

$$\rightarrow P(X=k) = p(x)$$

$$\cdot \sum_x p(x) = 1$$



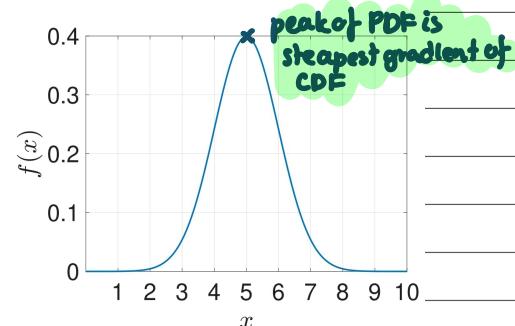
#### Probability Density Function (PDF)

- tell us the probability of CONTINUOUS random variable within a certain range.

$$f(x) = \frac{dF(x)}{dx}$$

$$\rightarrow P(a \leq X \leq b) = \int_a^b f(x) dx$$

$$\cdot \int_{-\infty}^{\infty} f(x) dx = 1$$

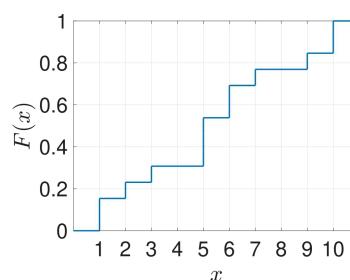


note that saying  $P(X=a)$  for PDF will mean nothing, we need a RANGE  
 $P(X=a) = 0$   
 $P(a \leq X \leq b) \neq 0$

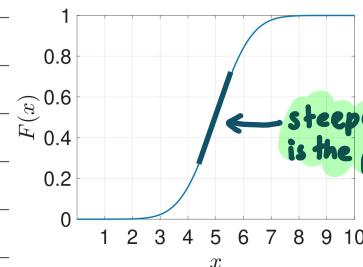
#### Cumulative Distribution Function (CDF)

$$\cdot F(k) = P(X \leq k)$$

for discrete data:



for continuous data:

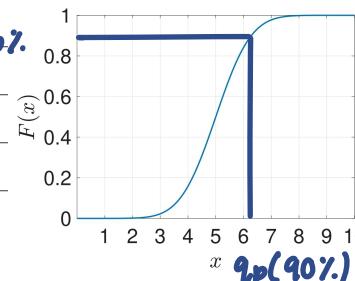


#### quantile function, $q_p(x)$

• quantile function is the inverse of CDF.

$$\begin{aligned} q_p(P/100) = u & \quad P(X \leq u) = P/100 \\ F(u) = P/100 & \end{aligned} \quad \left\{ q_p(x) = F^{-1}(x) \right. \quad \star$$

$p\% = 90\%$ .



↳ note that:

$$q_p(x) = F^{-1}(x) \text{ is exactly as } q_p(\frac{P}{100}) = F^{-1}(\frac{P}{100})$$

we usually write the second form as  $x = P/100$

$$q_p(90\%) = 6.25$$

eg.  $f(x) = \frac{\pi}{20} \sin\left(\frac{\pi x}{10}\right)$  for  $0 \leq x \leq 10$ , find 90<sup>th</sup> percentile.

Step 1: Find  $F(x)$

from  $0 \leq x \leq 10$ :

$$\begin{aligned} F(x) &= P(X \leq x) = \int_0^x \frac{\pi}{20} \sin\left(\frac{\pi t}{10}\right) dt \\ &= \frac{\pi}{20} \left[ -\frac{10}{\pi} \cos\left(\frac{\pi t}{10}\right) \right]_0^x \\ &= \frac{1}{2} - \frac{1}{2} \cos\left(\frac{\pi x}{10}\right), \quad x \in [0, 10] \end{aligned}$$

Step 2: find  $q_p(x)$

$$q_p(p/100) = F^{-1}(p/100)$$

$$\text{let } F^{-1}(x) = y$$

$$F(y) = x$$

$$\frac{1}{2} - \frac{1}{2} \cos\left(\frac{\pi y}{10}\right) = x$$

$$y = \frac{10}{\pi} \cos^{-1}(1-2x)$$

$$F^{-1}(x) = \frac{10}{\pi} \cos^{-1}(1-2x)$$

$$\therefore q_p(p/100) = \frac{10}{\pi} \cos^{-1}\left(1 - \frac{p}{100}\right)$$

Step 3: Sub  $p = 90$ :

$$q_{90}(90/100) = \frac{10}{\pi} \cos^{-1}\left(1 - \frac{90}{100}\right) = 7.952,$$

NOT in Y2!

## 3.2 Function of Random Variable (Conversion of PMF, PDF and CDF)

### 1. PMF conversion.

$$\text{let } Y = h(X)$$

$$P_Y(y) = P_X(h^{-1}(y))$$

$$\left\{ \begin{array}{l} P_Y(y) = P(Y=y) \\ = P(h(X)=y) \\ = P(X=h^{-1}(y)) \\ = P_X(h^{-1}(y)) \end{array} \right.$$

eg. if we have:

$$P_X(x) = \begin{cases} x/3, & x=0 \\ 1-x/3, & x=1 \end{cases}$$

and  $Y = 2X$

$$P_Y(y) = P_X(Y=y)$$

$$= P_X(2X=y)$$

$$= P_X(X=y/2)$$

$$= \begin{cases} y/6, & y=0 \\ 1-y/6, & y=2 \end{cases}$$

### 2. PDF conversion.

$$\text{let } Y = h(X)$$

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right|^{-1}$$

or usually written as.

$$f_Y(y) |dy| = f_X(x) |dx|$$

(absolute values are dropped, but remember  $dx$   $dy$  considered positive.)

### 3. CDF conversion

$$\text{let } Y = h(X)$$

$$F_Y(y) = P(Y \leq y)$$

$$= P(h(X) \leq y)$$

~~~~~ depends on  $h(X)$  to continue.  
as there is an inequality.

$$\text{eg. } h(X) = X+2$$

$$F_Y(y) = P(X+2 \leq y) = P(X \leq y-2) = F_X(y-2)$$

$$\text{eg. } h(X) = -X+3$$

$$F_Y(y) = P(-X+3 \leq y) = P(X \geq 3-y) = 1 - P(X < 3-y) = 1 - F_X(3-y)$$

not in syllabus,  
included  
to complete

## 3.3 Moments

## I. Moment (order of m)

replace  
X  
with  
 $X - E(X)$

$$E(X^m) = \sum_x x^m p(x) \quad \text{or} \quad E(X^m) = \int_{-\infty}^{\infty} x^m f(x) dx$$

(discrete)
(continuous)

## 2. Centred Moment (order of $m$ )

$$E[(X-E(X))^m] = \sum_{\text{(discrete)}}^{\infty} (x-E(x))^m p(x) \quad \text{or} \quad E[(X-E(X))^m] = \int_{\text{(continuous)}}^{-\infty} (x-E(x))^m f(x) dx$$

1st moment is actually just mean:

$$E(x') = E(x)$$

2<sup>nd</sup> centred moment is actually just variance:

$$E[(X-E(X))^2] = \text{Var}(X)$$

$\xrightarrow{\text{simplified to}} E(X^2) - (E(X))^2$

3<sup>rd</sup> centred moment, divided by  $s^3$  is the skewness.

$$g_1(x) = \frac{E[(x - E(x))^3]}{(s(x))^3}$$

$$\text{eg. } E[(h(x) - E(h(x)))^3]$$

**Moments of  $Y = h(x)$**  (just sub  $Y = h(x)$ !)  $\rightarrow$  same with central moment.

$$E(Y^m) = E(h(x)^m) = \sum_{x} h(x)^m p(x) \quad \text{or} \quad E(Y^m) = E(h(x)^m) = \int_{-\infty}^{\infty} h(x)^m f(x) dx$$

(discrete)      PMF of X      (continuous)      PDF of X

Y = h(x)      Y = h(x)      not the same

## 3.4 Expectation Algebra.

$$E(a) = a$$

$$Var(a) = 0$$

$$E(aX) = aE(X)$$

$$Var(aX) = a^2 Var(X)$$

$$E(a \pm X) = a \pm E(X)$$

$$Var(a \pm X) = Var(X)$$

$$E(a \pm bX) = a \pm bE(X)$$

$$Var(a \pm bX) = b^2 Var(X)$$

$$E(ax+by+c) = aE(X) + bE(Y) + c$$

$$\text{Var}(ax + by + c) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

if independent!

## C4. Probability Distribution.

### 4.1 Discrete Distribution.

1. Bernoulli Distribution. (1 trial, 2 outcomes)

$$X \sim B(1, p)$$

$$P(X=k) = p^k (1-p)^{1-k} \text{ for } k \in \{0, 1\}$$

} since  $n=1$ ,  $k$  can only be 0 or 1  
(either happen or no)

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

increase no.  
of trials

2. Binomial Distribution (n trial, each trial INDEPENDENT, 2 outcomes)

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

$$P(X=k) = {}^n C_k p^k (1-p)^{n-k}, \text{ for } k \in \{0, 1, \dots, n\}$$

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

as  $n$   
approaches  
infinity.

3. Poisson Distribution. (number of INDEPENDENT occurrence over a reference time)

$$X \sim Po(\mu)$$

( $\mu$  is the average number of occurrence over a reference time)

$$P(X=k) = \frac{e^{-\mu} \mu^k}{k!}, \text{ for } k \in \{0, 1, \dots\}$$

$$E(X) = \mu, \text{Var}(X) = \mu$$

4. Geometric Distribution. (infinite trial, each trial INDEPENDENT, 2 outcomes only, sequences of failure followed by one success)

$$X \sim Geo(p)$$

$P(X=k) = p(1-p)^k$   
(if  $k$  is number of failure  
till success)

or  $P(X=k) = p(1-p)^{k-1}$  for  $k \in \{0, 1, \dots\}$

(if  $k$  is number of trials till success)

$$P(X > k) = (1-p)^k$$

$$E(X) = \frac{1-p}{p} ; \text{Var}(X) = \frac{1-p}{p^2}$$

$$E(X) = \frac{1}{p} ; \text{Var}(X) = \frac{1-p}{p^2}$$

Both are correct!

This is the version given in data sheets.

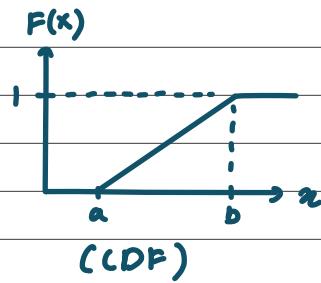
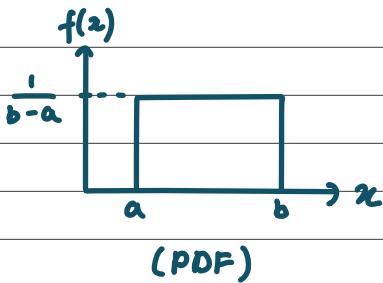
## 4.2 Continuous Distribution.

1. Uniform Distribution. (same probability over a given interval)

$$X \sim U[a, b]$$

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{for } x \in [a, b] \\ 0, & \text{otherwise.} \end{cases}$$

$$E(X) = \frac{a+b}{2}; \text{Var}(X) = \frac{(b-a)^2}{12}$$



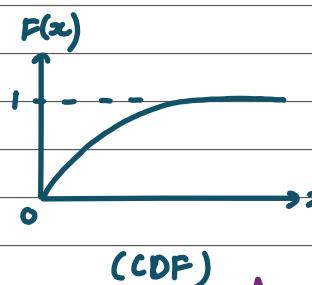
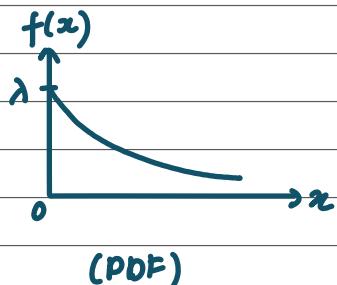
★ Similar to geometric

2. Exponential Distribution. (time taken till INDEPENDENT events occur at rate  $\lambda$ ,  $\lambda$  is no. of occurrence per unit time)

$$X \sim Exp(\lambda)$$

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

$$E(X) = \frac{1}{\lambda}; \text{Var}(X) = \frac{1}{\lambda^2}$$



$$X \sim Po(\mu), \quad \mu \text{ is no. of occurrence per reference time} \\ P(X=0) = \frac{e^{-\lambda x} (\lambda x)^0}{0!} = e^{-\lambda x} \quad \leftarrow \text{scale } \lambda \text{ to } \mu$$

$$\left\{ \begin{array}{l} X \sim Exp(\lambda) \\ P(X>x) = 1 - F(x) \\ \text{since } F(x) = \int \lambda e^{-\lambda x} dx \\ \therefore = e^{-\lambda x} \end{array} \right. \quad \text{same.}$$

$$X \sim Po(\lambda t), \quad Y \sim Exp(\lambda) \\ P(X=0) = P(Y>t)$$

★  $X \sim Po(\mu)$  and  $X \sim Exp(\lambda)$

$\mu$  and  $\lambda$  difference:  
eg  $\mu$ : study three times in 8 days.  
 $\lambda$ : study at rate  $0.375 \text{ day}^{-1}$

3. Gamma Distribution.

$$X \sim Gamma(d, \beta) \quad (d \text{ and } \beta \text{ are shape and scale parameters})$$

$$f(x) = \begin{cases} \frac{x^{d-1} e^{-x/\beta}}{\beta^d \Gamma(d)}, & \text{for } x > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\star \text{when } d=1, X \sim Gamma(1, \beta) = X \sim Exp(\frac{1}{\beta})$$

$$E(X) = d\beta; \text{Var}(X) = d\beta^2$$

$$\Gamma(d) = \int_0^\infty z^{d-1} e^{-z} dz, d \in \mathbb{R}$$

for  $d \in \mathbb{Z} : \Gamma(d) = (d-1)!$

$\left. \begin{array}{l} \text{basically } \Gamma(d) \text{ function extends} \\ \text{factorial to real domain.} \end{array} \right\}$

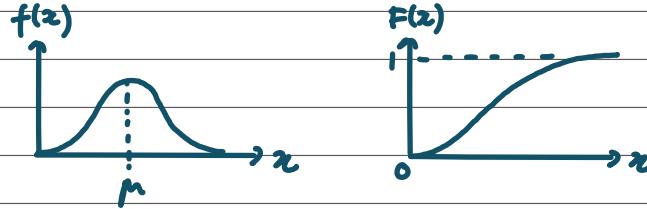
Gamma distribution graph varies a lot as  $d$  and  $\beta$  determine its shape.  
 for example, when  $d=1, \beta=1/\lambda, X \sim \text{Gamma}(1, 1/\lambda) = X \sim \text{Exp}(\lambda) !$   
 $X \sim \text{Gamma}(k/2, 2) = X \sim \chi^2(k) !$

#### 4. Normal Distribution (symmetry about $\mu$ , deviation with $\sigma$ among data)

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

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \forall x \in \mathbb{R}$$

$$E(X) = \mu; \text{Var}(X) = \sigma^2$$



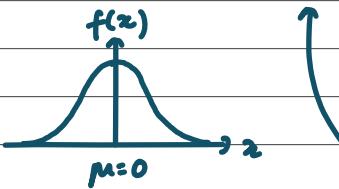
However,  $f(x)$  for  $X \sim N(\mu, \sigma^2)$  is complicated to compute so:  
**Standard Normal Distribution.**

from  $X \sim N(\mu, \sigma^2)$

$$\text{let } z = \frac{x-\mu}{\sigma}$$

$$\text{to } z \sim N(0,1) \text{ or } \left(\frac{x-\mu}{\sigma}\right) \sim N(0,1)$$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{z^2}{2}}, \forall z \in \mathbb{R}$$



still hard to compute?  
 the CDF,  $F(z)$  or known as  
 $\Phi(z)$  can be referred from  
 'Normal Distribution Table':

NORMAL CUMULATIVE DISTRIBUTION FUNCTION

| $x$ | 0.00   | 0.01   | 0.02   | 0.03   | 0.04   | 0.05   | 0.06   | 0.07   | 0.08   | 0.09   |
|-----|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.0 | 0.5000 | 0.5040 | 0.5120 | 0.5160 | 0.5199 | 0.5239 | 0.5279 | 0.5319 | 0.5359 |        |
| 0.1 | 0.5398 | 0.5438 | 0.5478 | 0.5517 | 0.5557 | 0.5596 | 0.5636 | 0.5675 | 0.5714 | 0.5753 |
| 0.2 | 0.5793 | 0.5832 | 0.5871 | 0.5910 | 0.5948 | 0.5987 | 0.6026 | 0.6064 | 0.6103 | 0.6141 |
| 0.3 | 0.6179 | 0.6217 | 0.6255 | 0.6293 | 0.6331 | 0.6368 | 0.6406 | 0.6443 | 0.6480 | 0.6517 |
| 0.4 | 0.6554 | 0.6591 | 0.6628 | 0.6664 | 0.6700 | 0.6736 | 0.6772 | 0.6808 | 0.6844 | 0.6879 |
| 0.5 | 0.6915 | 0.6950 | 0.6985 | 0.7019 | 0.7054 | 0.7088 | 0.7123 | 0.7157 | 0.7190 | 0.7224 |
| 0.6 | 0.7257 | 0.7291 | 0.7324 | 0.7357 | 0.7390 | 0.7422 | 0.7454 | 0.7486 | 0.7517 | 0.7549 |
| 0.7 | 0.7580 | 0.7611 | 0.7642 | 0.7673 | 0.7703 | 0.7734 | 0.7764 | 0.7794 | 0.7823 | 0.7852 |
| 0.8 | 0.7881 | 0.7910 | 0.7939 | 0.7967 | 0.7994 | 0.8023 | 0.8051 | 0.8078 | 0.8106 | 0.8133 |
| 0.9 | 0.8159 | 0.8186 | 0.8212 | 0.8238 | 0.8264 | 0.8289 | 0.8315 | 0.8343 | 0.8365 | 0.8389 |
| 1.0 | 0.8413 | 0.8438 | 0.8461 | 0.8485 | 0.8508 | 0.8531 | 0.8554 | 0.8577 | 0.8599 | 0.8621 |
| 1.1 | 0.8643 | 0.8665 | 0.8686 | 0.8708 | 0.8729 | 0.8749 | 0.8770 | 0.8790 | 0.8810 | 0.8830 |
| 1.2 | 0.8849 | 0.8869 | 0.8888 | 0.8907 | 0.8925 | 0.8944 | 0.8962 | 0.8980 | 0.8997 | 0.9015 |
| 1.3 | 0.9032 | 0.9049 | 0.9066 | 0.9082 | 0.9099 | 0.9115 | 0.9131 | 0.9147 | 0.9162 | 0.9177 |
| 1.4 | 0.9192 | 0.9207 | 0.9222 | 0.9236 | 0.9251 | 0.9265 | 0.9279 | 0.9292 | 0.9306 | 0.9319 |
| 1.5 | 0.9332 | 0.9345 | 0.9357 | 0.9370 | 0.9382 | 0.9394 | 0.9406 | 0.9418 | 0.9429 | 0.9441 |
| 1.6 | 0.9452 | 0.9463 | 0.9474 | 0.9484 | 0.9495 | 0.9505 | 0.9515 | 0.9525 | 0.9535 | 0.9545 |

e.g.  $\Phi(0.56) = 0.7123$ .

## eg1. Exponential Distribution Question (can be done with Poisson !)

During the period 1836 – 1961, 16 earthquakes of magnitude 6 or more on the Richter scale happened in San Francisco. What is the probability that an earthquake of magnitude 6 or more does not occur in the next 10 years?

### Exponential Distribution

$$X \sim \text{Exp}(\lambda)$$

$\lambda$  is the rate of occurrence ie:  
(number of occurrence per unit time)

$$\lambda = \frac{16}{1961 - 1836} = 0.128 \text{ per year}$$

$$P(X > 10) = 1 - P(X \leq 10) = 1 - \underbrace{F(10)}_{\text{CDF}}$$

$$f(z) = \begin{cases} \lambda e^{-\lambda z}, & z \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

for  $z \geq 0$ ,

$$F(z) = \int_0^z \lambda e^{-\lambda t} dt = [-e^{-\lambda t}]_0^z = 1 - e^{-\lambda z}$$

$$\begin{aligned} P(X > 10) &= 1 - F(10) \\ &= 1 - (1 - e^{-0.128(10)}) \\ &= 0.278 \end{aligned}$$

### Poisson Distribution.

$$X \sim \text{Po}(\mu)$$

$\mu$  is number of occurrence per reference time.

$$\mu = 0.128 \times 10 = 1.28 \text{ per ten year.}$$

$$P(X=0) = \frac{e^{-\mu} \mu^0}{0!}$$

$$\begin{aligned} &\text{'X in Poisson} \\ &\text{is } \mu \text{ occurrence} \\ &\text{in reference} \\ &\text{time'} \\ &= \frac{e^{-1.28} 1.28^0}{0!} \end{aligned}$$

$$X=0 \text{ means } 0 \text{ occurrence}$$

$\text{Earthquake in}$   
 $\text{reference time}$   
 $(10 \text{ years})$

## 4.3 Derivation of $E(X)$ and $\text{Var}(X)$ for all distribution

### 4.3.1 Discrete Distribution.

$$E(X) = \sum_{i \in S} x_i P(X=x_i), \quad \text{Var}(X) = \sum_{i \in S} x_i^2 P(X=x_i) - [ \sum_{i \in S} x_i P(X=x_i) ]^2$$

← ALWAYS KNOW YOUR DISTRIBUTION DOMAIN!

#### 1. Bernoulli Distribution.

$$X \sim B(1, p)$$

$$P(X=k) = p^k (1-p)^{1-k} \text{ for } k \in \{0, 1\}$$

$$\begin{aligned} E(X) &= \sum_{k=0}^1 k P(X=k) \\ &\text{since } k \in \{0, 1\}, \\ E(X) &= 0 \cdot P(X=0) + 1 \cdot P(X=1) \\ &\rightarrow 0 \cdot p^0 (1-p)^{1-0} + 1 \cdot p^1 (1-p)^{1-1} \\ &= p \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= \sum_{k=0}^1 k^2 P(X=k) - E(X)^2 \\ &= 0^2 P(X=0) + 1^2 P(X=1) - p^2 \\ &= p - p^2 \\ &= p(1-p) \end{aligned}$$

#### 2. Binomial Distribution

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

$$P(X=k) = {}^n C_k p^k (1-p)^{n-k}, \quad \text{for } k \in \{0, 1, \dots, n\}$$

$$\begin{aligned} E(X) &= \sum_{k=0}^n k P(X=k) \\ &\text{but } k=0 \text{ doesn't contribute ( } 0 \times \text{anything} = 0 \text{), so why not start at 1.} \end{aligned}$$

$$\begin{aligned} &= \sum_{k=1}^n k {}^n C_k p^k (1-p)^{n-k} \\ &= \sum_{k=1}^n k \frac{n!}{(n-k)! k!} p^k (1-p)^{n-k} \quad \leftarrow k! = k(k-1)!! \\ &= \sum_{k=1}^n \frac{n!}{(n-k)! (k-1)!!} p^k (1-p)^{n-k} \\ &\text{look a bit like } {}^{n-1} C_{k-1} = \frac{(n-1)!}{(n-1-k+1)! (k-1)!} = \frac{(n-1)!}{(n-k)! (k-1)!} \\ &= \sum_{k=1}^n n {}^{n-1} C_{k-1} p^k (1-p)^{n-k} \quad \leftarrow \text{since till this step, everything got } (-1) \\ &\quad \text{why not } p^k = p \cdot p^{k-1} ? \end{aligned}$$

$$\begin{aligned}
 &= np \sum_{k=1}^n {}^{n-1}C_{k-1} p^{k-1} (1-p)^{n-k} \\
 &= np \sum_{k=1}^n {}^{n-1}C_{k-1} p^{k-1} (1-p)^{(n-1)-(k-1)} \\
 &\stackrel{k=n}{\sum} \stackrel{k-1=n-1}{=} \sum_{k=0}^{k-1=n-1} {}^{n-1}C_{k-1} p^{k-1} (1-p)^{(n-1)-(k-1)} \\
 &= np \sum_{k=0}^{k-1=n-1} {}^{n-1}C_{k-1} p^{k-1} (1-p)^{(n-1)-(k-1)} \\
 &\quad \text{binomial formula} \rightarrow [p + (1-p)]^{n-1} \\
 &= 1^{n-1} \\
 &= 1
 \end{aligned}$$

$$\text{Var}(X) = E(X^2) - E(X)^2$$

$$= \sum_{k=0}^n k^2 {}^nC_k p^k (1-p)^{n-k} - (np)^2$$

↙ same steps as before.

$$= np \sum_{k=1}^n k {}^{n-1}C_{k-1} p^{k-1} (1-p)^{(n-1)-(k-1)}$$

this is a little bit hard...

### 3. Poisson Distribution.

$$X \sim Po(\mu)$$

$$P(X=k) = \frac{e^{-\mu} \mu^k}{k!}, \text{ for } k \geq 0$$

constant.

$$\begin{aligned} E(X) &= \sum_{k=0}^{\infty} k \frac{e^{-\mu} \mu^k}{k!} \\ &= e^{-\mu} \sum_{k=1}^{\infty} \frac{\mu^k}{(k-1)!} \quad \text{let } n=k-1 \\ &= \mu e^{-\mu} \sum_{n=0}^{\infty} \frac{\mu^{n+1}}{n!} \\ &= \mu e^{-\mu} e^{\mu} \\ &= \mu \end{aligned}$$

$k=0$  doesn't contribute.

through Taylor's series of  $e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$

$$\begin{aligned} \text{Var}(X) &= \sum_{k=0}^{\infty} k^2 \frac{e^{-\mu} \mu^k}{k!} - \mu^2 \\ &= \mu e^{-\mu} \sum_{k=1}^{\infty} k \frac{\mu^{k-1}}{(k-1)!} \quad k = k-1+1 \\ &= \mu e^{-\mu} \left( \sum_{k=2}^{\infty} (k-1) \frac{\mu^{k-1}}{(k-1)!} + \sum_{k=1}^{\infty} \frac{\mu^{k-1}}{(k-1)!} \right) \\ &\quad - \mu^2 \\ &\quad \text{when } k=1, \text{ the whole thing is zero.} \\ &\quad \text{so start from 2.} \end{aligned}$$

$\mu^{k-1} = \mu^{k-2} \mu$

$\sum_{k=2}^{\infty} (k-1) \frac{\mu^{k-1}}{(k-1)!} + \sum_{k=1}^{\infty} \frac{\mu^{k-1}}{(k-1)!}$

$= \mu^2 e^{-\mu} + \mu e^{-\mu} - \mu^2$

$= \mu /$

### 4. Geometric Distribution.

$$X \sim Geo(p)$$

$$P(X=k) = p(1-p)^{k-1} \quad \text{for } k \in \{0, 1, \dots\} \quad (\text{k is number of trials till success})$$

$$E(X) = \sum_{k=0}^{\infty} k p(1-p)^{k-1}$$

$$= p \sum_{k=1}^{\infty} k(1-p)^{k-1} \quad \text{looks like: } \frac{d}{dp} [-(1-p)^k]$$

$$= p \sum_{k=1}^{\infty} \frac{d}{dp} [-(1-p)^k]$$

$$\begin{aligned}
 &= -p \frac{d}{dp} \sum_{k=1}^{\infty} (1-p)^k \\
 &= -p \frac{d}{dp} \left( \frac{1-p}{p} \right) \\
 &= -p \left( -\frac{1}{p^2} \right) \\
 &= \frac{1}{p} //
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(X) &= \sum_{k=0}^{\infty} k^2 p (1-p)^{k-1} - \frac{1}{p^2} \\
 &= p \sum_{k=1}^{\infty} k^2 (1-p)^{k-1} \quad \text{it is like } \frac{d}{dp} [-k(1-p)^k] \\
 &= -p \frac{d}{dp} \sum_{k=1}^{\infty} k (1-p)^k - \frac{1}{p^2} \quad \text{this looks like } E(X) = \sum k p (1-p)^{k-1} \\
 &= -p \frac{d}{dp} \sum_{k=1}^{\infty} kp (1-p)^{k-1} \underbrace{\frac{(1-p)}{p}}_{\text{bring it out.}} - \frac{1}{p^2} \\
 &= -p \frac{d}{dp} \left( \frac{1-p}{p} \left( \frac{1}{p} \right) \right) - \frac{1}{p^2} \\
 &= -p \frac{d}{dp} \left( \frac{1-p}{p^2} \right) - \frac{1}{p^2} \\
 &= -p \frac{d}{dp} (p^{-2} - p^{-1}) - p^{-2} \\
 &= -p(-2p^{-3} + p^{-2}) - p^{-2} \\
 &= 2p^{-2} - p^{-1} - p^{-2} \\
 &= p^{-2} - p^{-1} \\
 &= \frac{1}{p^2} - \frac{1}{p} \\
 &= \frac{1-p}{p^2} //
 \end{aligned}$$

#### 4.3.2 Continuous Distribution.

$$E(X) = \int_{\Omega} x f(x) dx, \quad \text{Var}(X) = \int_{\Omega} x^2 f(x) dx - \left[ \int_{\Omega} x f(x) dx \right]^2$$

## C5. Fitting Probability Distribution

### 5.1 Method of Moment (1 out of 2 method covered in yr2)

Step 1. Think which distribution we want to fit into our sample.  $\star$  (usually told in exam)

e.g.  $X \sim B(n, p)$ ,  $X \sim N(\mu, \sigma^2)$ , ...

Step 2. Equate population moment (our fitted distribution's) to the unbiased sample moment (the original data sets)

$$\left. \begin{array}{ll} E(X) = \bar{x} & (\text{moment 1}) \\ \text{Var}(X) = s_x^2 & (\text{centred moment 2}) \\ g_1(X) = \hat{g}_1(x) & (\text{centred moment 3}) \\ \vdots & \vdots \\ \text{population} & \text{sample (have to find)} \\ (\text{have to memorise}) & \end{array} \right\}$$

$\star$  IF NUMBER OF EQUATION AVAILABLE IS GREATER THAN THE NUMBER OF PARAMETER OF OUR FITTED DISTRIBUTION, WE NEED TO DO MINIMISATION WITH WEIGHTAGE INSTEAD

$\star$  if in exams, was not told how many eqn of moment required, choose exactly the same:  
 $X \sim B(n, p)$ :  $E(X) = np = \bar{x}$   
 2 parameters       $\text{Var}(X) = np(1-p) = s_x^2$   
 2 eqns

Step 3. Minimisation process with weightage.

$$\min \sum_{k=1}^p w_k (M_k - m_k)^2 = \min \left\{ w_1 (E(X) - \bar{x})^2 + w_2 (\text{Var}(X) - s_x^2)^2 + \dots \right\}$$

$\ominus$  population moments       $\uparrow$  sample moments

$\star$   $w$  (weightage) are usually given in exams.  
 But if not, smaller moment should have higher weightage.

e.g. if there is  $w_1, w_2$  and  $w_3$ :

$$w_1 = 0.5, w_2 = 0.3, w_3 = 0.2$$

$\star$  Why when no of moment eqn = no of parameters is enough but sometimes we use more no of moment eqn? (which is a lot harder to do)  
 cause the more no of moment eqn we use, the more accurate the fit is!

## eg1. (let us use no of moment eqns = no of parameters)

Using the method of moments, fit a uniform distribution to the sample below

$$\{1.0, 1.1, 1.2, 1.1, 1.3, 1.0, 1.7, 2.0\}$$

### Step 1. Determine the fitted probability distribution

$$X \sim U[a, b] \quad \leftarrow 2 \text{ parameters}$$

$$1^{\text{st}} \text{ moment : } E(X) = \frac{a+b}{2}$$

$$2^{\text{nd}} \text{ moment : } \text{Var}(X) = \frac{(b-a)^2}{12}$$

### Step 2. Equate population moment to sample moment.

$$\bar{x} = \frac{1.0 + 1.1 + 1.2 + 1.1 + 1.3 + 1.0 + 1.7 + 2.0}{8} = 1.3 \quad \left. \right\} 2 \text{ eqns of moment}$$

$$(\text{unbiased}) S_x^2 = 0.3625$$

$$\frac{a+b}{2} = 1.3 \quad \frac{(b-a)^2}{12} = 0.3625$$

$$\therefore a = 0.2572, b = 2.3428$$

$$X \sim U[0.2572, 2.3428]$$

## eg2. (Let us try no of moment eqn > no of parameters)

Using the method of moments, fit a Poisson distribution to the sample below of numbers of monthly nuclear alerts at a power plant:

$$\{0, 1, 1, 0, 0, 2, 0, 1\}$$

### Step 1. Determine the fitted probability distribution

$$X \sim Po(\mu)$$

$$1^{\text{st}} \text{ moment : } E(X) = \mu$$

$$2^{\text{nd}} \text{ moment : } \text{Var}(X) = \mu$$

$$\bar{x} = \frac{5}{8}, S_x^2 = 0.5536$$

WE SKIP STEP TWO CAUSE:

NO OF MOMENT EQN > NO OF PARAMETERS!

Step 3. Minimisation process with weightage.  
(for this question, we use  $w_1 = 1, w_2 = 1$ )

$$\text{let } M(\mu) = 1\left(\mu - \frac{5}{8}\right)^2 + 1\left(\mu - 0.5536\right)^2$$

$$\text{to find min of } M: \frac{dM}{d\mu} = 0 \left( \frac{d^2M}{d\mu^2} > 0 \right)$$

$$\frac{dM}{d\mu} = 2\left(\mu - \frac{5}{8}\right) + 2\left(\mu - 0.5536\right)$$

$$0 = 2\left(\mu - \frac{5}{8}\right) + 2\left(\mu - 0.5536\right)$$

$$\therefore \mu = 0.585 \quad X \sim Po(0.585)$$

$$\left( \begin{array}{l} \text{just to check } \mu = 0.585 \text{ will result in } M \text{ as minimum and not maximum:} \\ \frac{d^2M}{d\mu^2} = 2+2=4>0 \end{array} \right)$$

eg3. (Let us use 3 eqn of moments) ☒ IGNORE THIS QUESTION  
Fit, with method of moment: (binomial) not in year 2 stats!

$$[2, 1, 3, 2, 4, 0, 2, 3, 1, 2]$$

$$E(X) = np \quad \bar{x} = 2$$

$$\text{Var}(X) = np(1-p) \quad S_{2x}^2 = \frac{12}{9}$$

$$g_1(x) = \frac{1-2p}{\sqrt{np(1-p)}} \quad \hat{g}_1(x) = 0$$

$$\text{let } S(n, p) = w_1(np - \bar{x})^2 + w_2(np(1-p) - S^2)^2 + w_3\left(\frac{1-2p}{\sqrt{np(1-p)}} - \hat{g}_1(x)\right)^2$$

let  $w_1, w_2, w_3 = 1$  (for simplicity)

to minimise a multivariable function, it is hard (partial) derivative

(alt way: test every possible n and p):

$\therefore n=5, p=0.4$  will result in minimal  $S_{\text{min}}$

$\therefore X \sim B(5, 0.4)$

## 5.2 Maximum Likelihood Method (2 out of 2 method covered in y2)

basically just:  $\underset{\Theta}{\text{Max}} \left\{ L(\Theta) = \prod_{i=1}^m f_X(x_i | \Theta) \right\}$  (continuous)

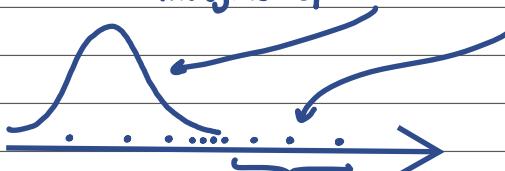
$\underset{\Theta}{\text{Max}} \left\{ L(\Theta) = \prod_{i=1}^m p_X(x_i | \Theta) \right\}$  (discrete)

looks complicated?

Basically it means, we want to find the parameter of a fitted distribution (eg. for  $X \sim B(n, p)$ ,  $\Theta$  is  $n$  and  $p$ ) that will result a maximum  $L(n, p)$ :  
 $L(n, p) = p(X=x_1) \times p(X=x_2) \times p(X=x_3) \times \dots$

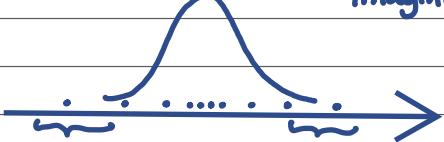
Why?

imagine we fit this PDF to this sample,



probability getting these dots are so small,  
when multiply all the probability you can feel that  
the  $L(\Theta)$  is very small!

imagine now we fit this PDF to this sample



although the probability on the sides are still small,  
but the majority is in the center, which results in higher probability!

So, how to find  $\Theta$  that correspond to  $L(\Theta)$  maximum?

if one parameter only, if more than one parameter,

$$\frac{dL(\Theta)}{d\Theta} = 0, \quad \frac{d^2L(\Theta)}{d\Theta^2} < 0 \quad \frac{\partial L(\Theta)}{\partial \Theta} = 0 \text{ for every } \Theta, \quad \frac{\partial^2 L(\Theta)}{\partial \Theta^2} < 0$$

(maxima)

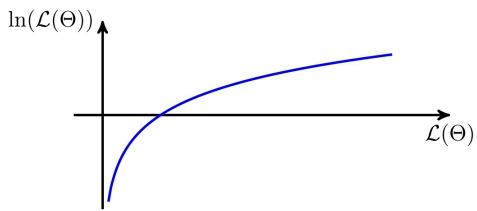
However,  $L(\Theta) = \prod_{i=1}^n f_X(x_i | \Theta)$  is hard to be differentiated (product rule too many products)

hence we find  $\underset{\Theta}{\text{Max}} \left\{ \ln L(\Theta) = \ln \left( \prod_{i=1}^n f_X(x_i | \Theta) \right) = \sum_{i=1}^n \ln(f_X(x_i | \Theta)) \right\}$

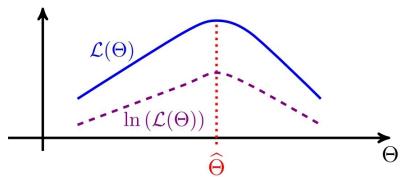
how does product become sum?  $\ln(ABC) = \ln A + \ln B + \ln C$

We can do that cause  $\ln(\mathcal{L}(\theta))$  and  $\mathcal{L}(\theta)$  shares the same maximum!

Log is a monotonically increasing function



This means that the log likelihood is maximum at the same point as the likelihood



eg 1.

(I didn't do  $\sum \ln(f_x(x_i | \theta))$  cause exponent are easy to multiply)

Using the maximum likelihood method, fit a Poisson distribution to the sample below of numbers of monthly nuclear alerts at a power plant:

$$\{0, 1, 1, 0, 0, 2, 0, 1\}$$

$$X \sim Po(\mu) \quad \underset{\approx}{\text{Max}} \left\{ \mathcal{L}(\theta) = \prod_{i=1}^m p_x(x_i | \theta) \right\} \quad (\text{discrete})$$

$$p_x(x_i | \mu) = \frac{\mu^{x_i} e^{-\mu}}{x_i!}, \quad \mathcal{L}(\mu) = \prod p_x(x_i | \mu)$$

inserting the data sets from sample,

$$\begin{aligned} \mathcal{L}(\mu) &= P(X=0)P(X=1)P(X=1)P(X=0)P(X=0)P(X=2)P(X=0)P(X=1) \\ &= P(X=0)^4 P(X=1)^3 P(X=2)^1 \\ &= \left( \frac{\mu^0 e^{-\mu}}{0!} \right)^4 \times \left( \frac{\mu^1 e^{-\mu}}{1!} \right)^3 \times \left( \frac{\mu^2 e^{-\mu}}{2!} \right)^1 \\ &= \frac{1}{2} \mu^5 e^{-8\mu} \end{aligned}$$

to find Max of  $\mathcal{L}(\mu)$ :

$$\frac{d\mathcal{L}(\mu)}{d\mu} = \frac{1}{2} (5\mu^4 e^{-8\mu} - 8\mu^5 e^{-8\mu}) \quad \left. \quad \frac{d^2\mathcal{L}(\mu)}{d\mu^2} < 0 \text{ (check passed)} \right\}$$

$$0 = \frac{1}{2} e^{-8\mu} (5\mu^4 - 8\mu^5)$$

$$\mu = \frac{5}{8}$$

$$\therefore X \sim Po\left(\frac{5}{8}\right)$$

2021q2c(i)

\* Usually max likelihood method in this course will only ask for distribution that has one unknown parameter:  
e.g.  $X \sim Exp(\lambda)$ ,  $X \sim Po(\mu)$ ,  $X \sim B(50, p)$

but there is an exception!

$X \sim U(a, b)$ , width =  $b-a$  if width = range of sample,  
 $L(w) = \prod \frac{1}{w}$   $L(w) \rightarrow \max!$

## 5.3 Goodness of Fit

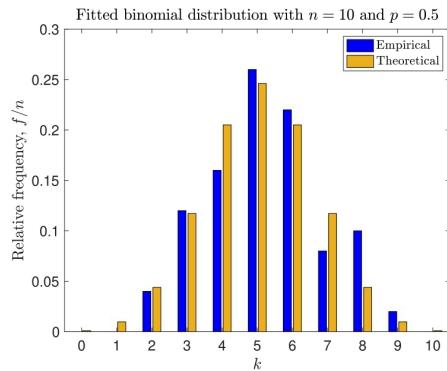
After fitting a distribution onto a sample, we should check how "good" of a fit the distribution is.

3 methods to do so:

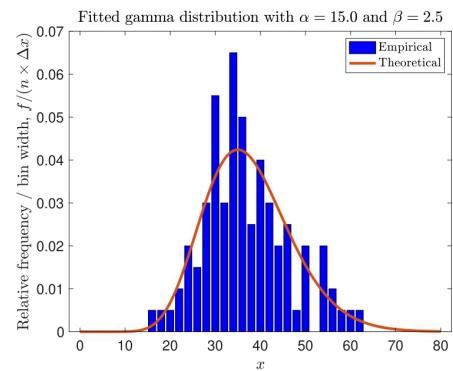
1. Compare sample's histogram to distribution's PDF (PMF if discrete)
2. Compare sample's CDF to distribution's CDF
3. Create a Quantile-Quantile Plot (Q-Q plot)

### 1. Compare sample's histogram to distribution's PDF (PMF if discrete)

Discrete  $X$ : use relative frequency 



Continuous  $X$ : use relative frequency divided by bin width 

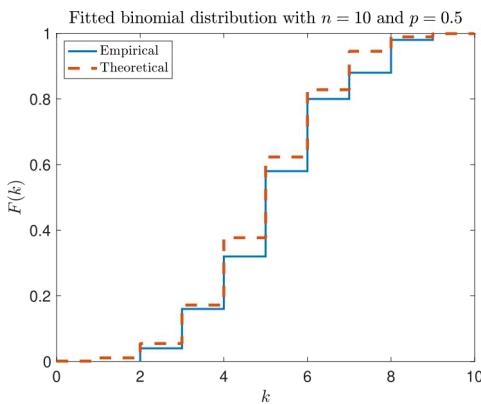


} empirical means  
Sample's, theoretical means  
fitted population distribution

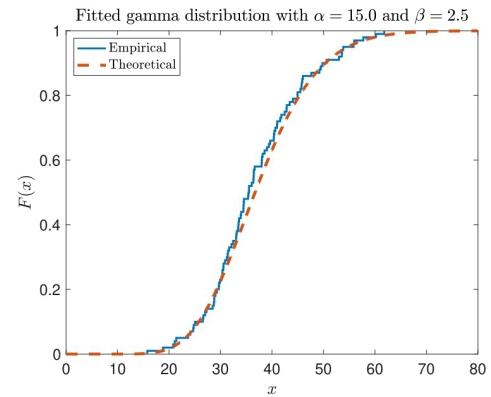
histogram with PMF (discrete) | histogram with PDF (continuous)

### 2. Compare sample's CDF to distribution's CDF

Discrete distribution



Continuous distribution



} empirical means  
Sample's, theoretical means  
fitted population distribution

### 3. Create a Quantile-Quantile Plot (QQ plot)

Having a QQ Plot let us able to compare quantile of the sample and the quantile of the fitted distribution (and check if they match)

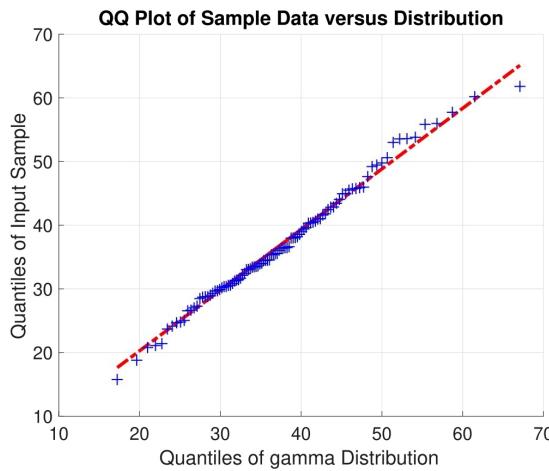
To plot QQ plot:

- vertical axis (y-coordinate) is just your sample, sorted from first to last.
- horizontal axis (x-coordinate), use:

$$q_P\left(\frac{P}{100}\right) = F^{-1}\left(\frac{P}{100}\right), \text{ where } \frac{P}{100} = \frac{j-0.5}{n}, \text{ for } j=1, 2, 3, \dots, n$$

so we need to find CDF (of the fitted distribution),  
then inverse it

Continuous  $X$ : Gamma distribution



} if it assembles a straight line,  
it is a good fit.

## c6. Multiple Variables

### 6.1 Covariance and Correlation

This chapter have a lot of info to memorise. Hence, I'll highlight all the part that should be memorise. (Other part that is not highlighted is the explanation)

#### 6.1.1 Covariance

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])] \neq \text{Cov}(X, X) = E[XX] - E[X]E[X]$$

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

wed to prove

$$= E[X^2] - E[X]^2$$

$$= \text{Var}(X)$$

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

$$\begin{aligned} &= \text{Cov}(X, X) + \text{Cov}(X, Y) + \text{Cov}(Y, X) + \text{Cov}(Y, Y) \\ &= \text{Var}(X) + 2\text{Cov}(X, Y) + \text{Var}(Y) \end{aligned} \quad \left. \begin{array}{l} \text{using the fact that:} \\ \text{Cov}(X, X) = \text{Var}(X) \end{array} \right.$$

1.2.4 Measure Relation between Datasets. (to check whether there is any LINEAR correlation between  $x_i$  and  $y_i$ )

**1. SAMPLE covariance.** (problem about this is covariance has units. if  $x$  and  $y$  has unit  $[m]$  then covariance has unit  $[m^2]$ )  
 $\text{cov}_{x,y} = \frac{1}{n-1} E[(x-\bar{x})(y-\bar{y})]$  so to non-dimensionalise...

**2. SAMPLE correlation coefficient.** bounded by:  $-1 \leq \text{cav}_{xy} \leq 1$

**if population:**

**1. Population covariance**

$$\text{Cov}(X, Y) = \frac{1}{n} E[(x-\mu_x)(y-\mu_y)] = E(XY) - E(X)E(Y)$$

**2. Population correlation coefficient.**

$$C_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

in c1.2.4, we looked at both sample covariance and population covariance.

In this chapter we only look at population covariance as  $X$  and  $Y$  are random variables.

#### 6.1.2 Properties of Covariance.

$$\begin{aligned} 1. \text{Cov}(cX, Y) &= c \text{Cov}(X, Y) \\ \text{Cov}(X, cY) &= c \text{Cov}(X, Y) \end{aligned} \quad \left. \begin{array}{l} \text{constant can be brought out!} \\ \text{Cov}(ax, bY) = ab \text{Cov}(X, Y) \end{array} \right.$$

$$2. \text{Cov}(X, Y) = \text{Cov}(Y, X) \quad \text{random variable can be swapped! (known as symmetric property)}$$

$$3. \text{Cov}(X, c) = 0 \quad \left. \begin{array}{l} \text{if any random variable is 'Cov' with a constant.} \\ \text{Cov}(c, Y) = 0 \quad \text{the Covariance is 0!} \end{array} \right.$$

$$4. \text{Cov}\left(\sum_{i=1}^n a_i X_i, \sum_{j=1}^m b_j Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m a_i b_j \text{Cov}(X_i, Y_j)$$

don't remember the formula!

looks complicated, but just remember:  $\text{Cov}(aX_1 + bX_2, cY_1 + dY_2)$   
 (do step by step)

$$\begin{aligned} &= \text{Cov}(aX_1, cY_1) + \text{Cov}(aX_1, dY_2) + \text{Cov}(bX_2, cY_1) + \text{Cov}(bX_2, dY_2) \\ &= ac \text{Cov}(X_1, Y_1) + ad \text{Cov}(X_1, Y_2) + bc \text{Cov}(X_2, Y_1) + bd \text{Cov}(X_2, Y_2) \end{aligned}$$

$\rightarrow$  COV is not the same as Cov.  
 Cov is Covariance MATRIX and Cov is just Covariance (a number)

### 5. $\text{COV}(X_1, X_2, X_3, \dots, X_n)$

not as important as the first 4.  
 (can skip)

$$= \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_n) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_n, X_1) & \text{Cov}(X_n, X_2) & \dots & \text{Var}(X_n) \end{pmatrix}$$

things to take note: 1.  $\text{Cov}(X_1, X_2) = \text{Cov}(X_2, X_1)$ , hence COV is symmetric about the main diagonal.

2. The main diagonal are all  $\text{Cov}(X, X) = \text{Var}(X)$ !

eg.

Let  $X$  and  $Y$  be two independent  $N(0, 1)$  random variables and

$$Z = 1 + X + XY^2, \\ W = 1 + X.$$

Find  $\text{Cov}(Z, W)$ .

$$\begin{aligned} \text{Cov}(Z, W) &= \text{Cov}(1 + X + XY^2, 1 + X) \\ &= \text{Cov}(1, 1) + \text{Cov}(1, X) + \text{Cov}(X, 1) + \text{Cov}(X, X) + \text{Cov}(XY^2, 1) + \text{Cov}(XY^2, X) \\ &\quad \{ \text{Cov}(X, X) = \text{Var}(X) \} \\ &= \text{Var}(X) + \text{Cov}(XY^2, X) \end{aligned}$$

$$\begin{aligned} &= \text{Var}(X) + E(XY^2) - E(XY^2)E(X) \\ &= \text{Var}(X) + E(X^2Y^2) - E(XY^2)E(X) \end{aligned}$$

after simplifying with all covariant property,  
 we use the covariance formula:

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

} now we finish simplifying into  
 only expectation and variance,  
 we can now use the fact that  $X \sim N(0, 1)$ ,  $Y \sim N(0, 1)$  ← from question.

$$\begin{aligned} E(X) &= 0 & E(Y) &= 0 \\ \text{Var}(X) &= 1 & \text{Var}(Y) &= 1 \end{aligned}$$

... and also independent means :

$$E[g(X)h(Y)] = E[g(X)]E[h(Y)] \leftarrow \text{will be covered later in c6.}$$

$$= \text{Var}(X) + E(X^2)E(Y^2) - E(X)E(Y^2)E(X)$$

$$\begin{aligned} &= 1 + 1 \times 1 + 0 \times 1 \times 0 \\ &= 2 \end{aligned}$$

$$\begin{aligned} \text{Var}(Y) &= E(Y^2) - E(Y)^2 \\ E(Y^2) &= \text{Var}(Y) + E(Y)^2 \\ &= 1 + 0^2 \\ &= 1 \\ E(X^2) &= \text{Var}(X) + E(X)^2 \\ &= 1 \end{aligned}$$

→ only measures LINEAR correlation!

### 6.1.3 Correlation Coefficient.

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

1.2.4 Measure Relation between Datasets. (to check whether there is any LINEAR correlation between  $x_1$  and  $y_1$ )

1. SAMPLE covariance. (problem about this is covariance has units. If  $x$  and  $y$  has unit [m] then covariance has unit [ $m^2$ ])

$$\text{cov}_{x,y} = \frac{1}{n-1} \sum (x - \bar{x})(y - \bar{y}) \quad \text{so to non-dimensionalise...}$$

2. SAMPLE correlation coefficient.

$$C_{x,y} = \frac{\text{cov}_{x,y}}{S_x S_y}, \quad -1 \leq C_{x,y} \leq 1$$

i.  $C_{x,y} = -1$ : (perfect negative linear correlation)

ii.  $C_{x,y} = 0$ : (completely no correlation)

iii.  $C_{x,y} = 1$ : (perfect positive linear correlation)

in c1.2.4, we looked at both sample correlation coefficient and population correlation coefficient. In this chapter we only look at population correlation coefft.  $X$  and  $Y$  are random variables.

### 6.1.4 Properties of Correlation Coefficient

1.  $-1 \leq \rho(X, Y) \leq 1$  (correlation coeff. is always between -1 and 1)

2.  $\rho(ax+b, cy+d) = \rho(X, Y)$  (multiply random variable with scalar or addition with constant has NO effect)

3. If  $\rho(X, Y) = 0$ , uncorrelated

If  $\rho(X, Y) > 0$ , positively correlated

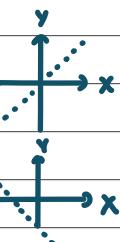
If  $\rho(X, Y) < 0$ , negatively correlated.

If  $\rho(X, Y) = 1$ ,  $Y = aX + b$

$$a > 0$$

If  $\rho(X, Y) = -1$ ,  $Y = aX + b$

$$a < 0$$



☆ uncorrelated is not same as independent!

if independent → uncorrelated ✓

but if uncorrelated, it doesn't mean it is independent!

(there's some special case like if both random variable are normally distributed)

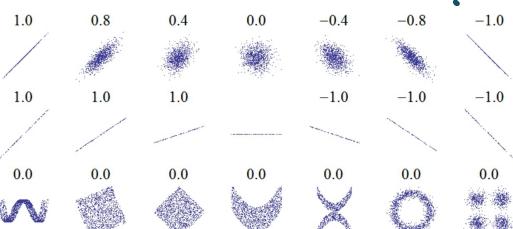
$\rho(X, Y) = 0$  means uncorrelated but doesn't mean independent!

any property independent have → uncorrelated have (but not the other way)

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

this is true if  $X$  and  $Y$  are independent, if independent → uncorrelated.

hence it is true if  $X$  and  $Y$  are uncorrelated!



note that it doesn't depend on slope

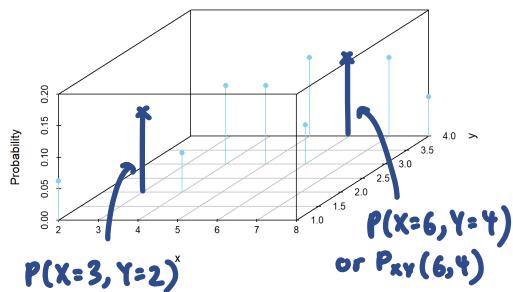
these are example of no correlation

## 6.2 Joint Distribution

It is just probability distribution in 2D! PMF for discrete usually given in table; PDF for continuous usually just given as a function  $f(x,y)$

### Joint Distribution

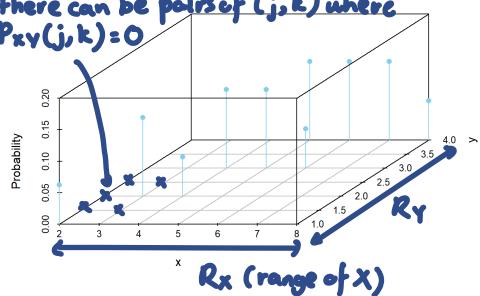
#### Discrete



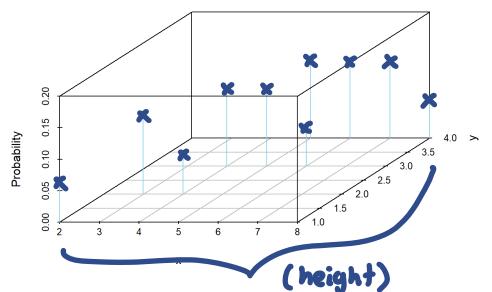
$$\text{Joint PMF: } P_{XY}(j,k) = P(X=j, Y=k)$$

no need to memorise {

there can be pairs of  $(j, k)$  where  $P_{XY}(j, k) = 0$



we usually define range as a rectangle,  $R_x \times R_y$  in size.

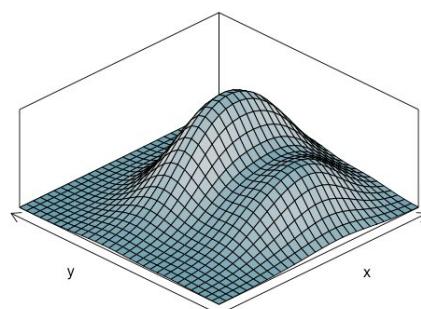


summing all the probability in the domains

$$\sum_{(j,k) \in R_{XY}} P_{XY}(j,k) = 1$$

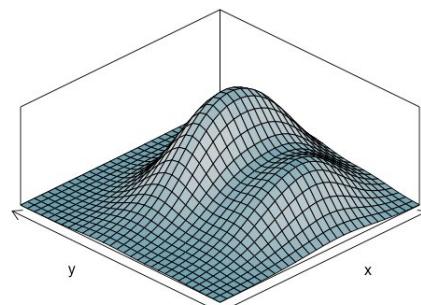
$$\sum_{k \in R_y} \sum_{j \in R_x} P_{XY}(j,k) = 1$$

#### Continuous

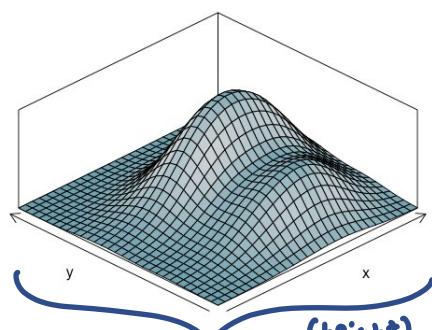


$$\text{Joint PDF: } f_{XY}(x,y) = \frac{P(z \leq X \leq z+dz, y \leq Y \leq y+dy)}{dz dy}$$

} no need to memorise.



continuous, range is from  $-\infty$  to  $\infty$  for both  $X$  and  $Y$ . (And can be trimmed into and shape, eg. Range:  $x^2 + y^2 < 9$ )



Summing all the probability

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x,y) dx dy = 1$$

## Joint Distribution

### Discrete

|         | $Y = 0$       | $Y = 1$       | $Y = 2$       |
|---------|---------------|---------------|---------------|
| $X = 0$ | $\frac{1}{6}$ | $\frac{1}{4}$ | $\frac{1}{8}$ |
| $X = 1$ | $\frac{1}{8}$ | $\frac{1}{6}$ | $\frac{1}{6}$ |

PMF usually expressed in a table

### Continuous

$$f_{XY}(x, y) = \begin{cases} c(x^2 + y^2) & \text{for } 0 \leq x \leq 1 \text{ and } 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

PDF usually expressed as a multivariable function  $f(x, y)$

for CDF:

$$F_{XY}(j, k) = P(X \leq j, Y \leq k)$$

for CDF: *remember CDF SS*

$$F_{XY}(x, y) = \int_{-\infty}^y \int_{-\infty}^x f_{XY}(u, v) du dv$$

this also means that differentiating CDF twice can get PDF too :

$$f_{XY}(x, y) = \frac{\partial^2}{\partial x \partial y} F_{XY}(x, y)$$

$u$  and  $v$  are just dummy variables

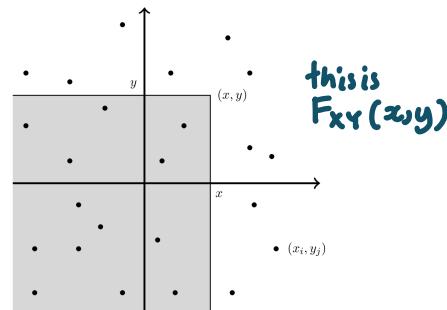
important results:

$$0 \leq F_{XY}(x, y) \leq 1$$

$$F_{XY}(\infty, \infty) = 1$$

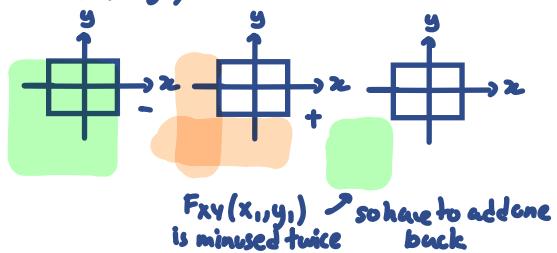
$$F_{XY}(-\infty, y) = 0 \text{ for any } y$$

$$F_{XY}(x, -\infty) = 0 \text{ for any } x$$



$$P(z_1 \leq X \leq z_2, y_1 \leq Y \leq y_2)$$

$$= F_{XY}(z_2, y_2) - F_{XY}(z_1, y_2) - F_{XY}(z_2, y_1) + F_{XY}(z_1, y_1)$$



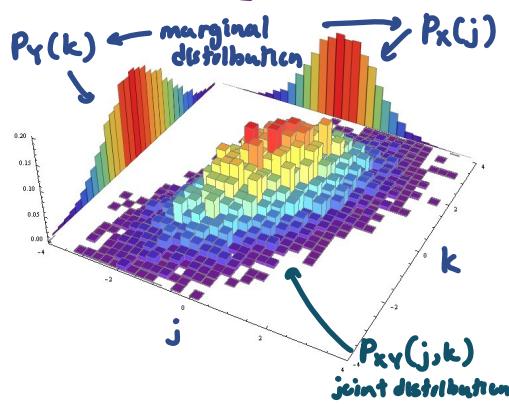
} DON'T  
MEMORISE  
understand it.

## 6.3 Marginal Distribution

### Marginal Distribution

basically it's when the distributions of the random variable considered separately

#### Discrete



$$P_x(j) = \sum_{k=0}^{\infty} P_{xy}(j, k)$$

for each  $j$ , sum all the  $k$   
e.g.  $P_x(1) = \sum_{k=0}^{\infty} P_{xy}(1, k)$

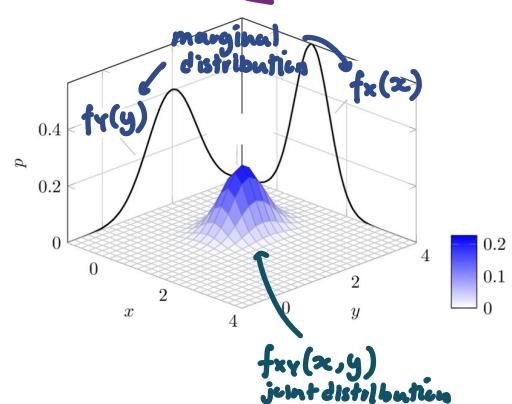
$$P_x(2) = \dots$$

⋮

all of these put in a table and you get the distribution.

$$P_y(j) = \sum_{j=0}^{\infty} P_{xy}(j, k)$$

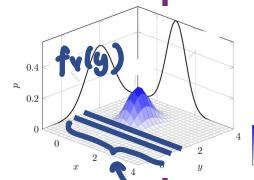
#### Continuous



$$f_x(x) = \int_{-\infty}^{\infty} f_{xy}(x, y) dy$$

similarly,

$$f_y(y) = \int_{-\infty}^{\infty} f_{xy}(x, y) dx$$



integrate all these small strips so  $\frac{dx}{2}$ .

#### ❖ FIRST WAY TO PROOF INDEPENDENT :

IF DISCRETE:

$$\underbrace{P_{xy}(j, k)}_{\text{joint distribution.}} = \underbrace{P_x(j) P_y(k)}_{\text{marginal dist.}}$$

IF CONTINUOUS:

$$\underbrace{f_{xy}(x, y)}_{\text{joint dist.}} = \underbrace{f_x(x) f_y(y)}_{\text{marginal dist.}}$$

instead of PDF, if continuous we can use CDF:

$$F_{xy}(x, y) = F_x(x) F_y(y)$$

## 6.4 Conditional Distribution.

| <u>Conditional Distribution</u>                                                                                                                                |                                                                 |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------|
| $P(A B) = \frac{P(A \cap B)}{P(B)}$ $\xrightarrow{\text{2 vars}}$ $P(X \in C   Y \in D) = \frac{P(X \in C, Y \in D)}{P(Y \in D)}$<br>C, D are range of X and Y |                                                                 |
| <u>Discrete</u><br>$P_{X Y}(j,k) = P(X=j   Y=k) = \frac{P(X=j, Y=k)}{P(Y=k)}$                                                                                  | <u>Continuous</u><br>$f_{XY}(x y) = \frac{f_{XY}(x,y)}{f_Y(y)}$ |
| $\therefore$ basically, Conditional = $\frac{\text{Joint}}{\text{Marginal}}$                                                                                   |                                                                 |

Q SECOND WAY TO PROOF INDEPENDENT :

IF DISCRETE :

$$P_{X|Y}(x|y) = P_X(x)$$

$$P_{Y|X}(y|x) = P_Y(y)$$

IF CONTINUOUS :

$$f_{X|Y}(x|y) = f_X(x)$$

$$f_{Y|X}(y|x) = f_Y(y)$$

recall that this is similar to :

$P(A|B) = P(A)$  proof that A and B are independent event.

eg. Find CDF of X given PDF of XY :  $f_{XY} = \begin{cases} xy, & 0 \leq y \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$

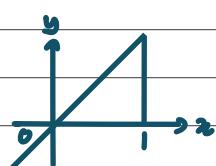
Step 1. Find PDF of X,  $f_X(x)$

$$f_X(x) = \int f_{XY}(x,y) dy$$

Bye RY

$$\begin{aligned} y &= x \\ &= \int x y dy \\ y_0 &= 0 \end{aligned}$$

$$= \frac{1}{2} x^3$$



Step 2: using  $f_X(x)$ , find  $F_X(x)$

$$F_X(x) = \int_0^x \left( \int_v^x \frac{1}{2} u^3 du \right) dv$$

$$= \int_0^x \left[ \frac{1}{8} u^4 \right]_v^x dv$$

$$= \frac{1}{8} \int_0^x x^4 - v^4 dv$$

$$= \frac{1}{8} \left[ x^5 - \frac{1}{5} v^5 \right]_0^x$$

$$= \frac{1}{8} \left( \frac{4}{5} x^5 \right) = \frac{1}{10} x^5$$

## 6.5 Important Identity

### 1. Conditional Expectation:

$$E[X|Y=y_j] = \sum_{x_i \in R_X} x_i P_{X|Y}(x_i|Y=y_j) \quad \text{from: } E(X) = \sum_{x_i \in R_X} x_i P(X=x_i)$$

$$X \rightarrow X|Y=y_j \quad x_i \rightarrow x_i|Y=y_j$$

### 2. Conditional Variance

$$\text{Var}(X|Y=y) = E[X^2|Y=y] - E[X|Y=y]^2 \quad \text{from: } \text{Var}(X) = E(X^2) - E(X)^2$$

$$x^2 \rightarrow x^2|Y=y$$

### 3. Law of Total Probability.

$$P_X(x) = \sum_{y_j \in R_Y} P_{XY}(x, y_j) = \sum_{y_j \in R_Y} P_{X|Y}(x|y_j) P_Y(y_j) \quad \text{from: } P(A) = \sum_{B_i} P(A \cap B_i) = \sum_{B_i} P(A|B_i) P(B_i)$$

$$A \rightarrow X=x \quad B_i \rightarrow Y=y_i$$

### 4. Law of Total Expectation.

$$E[X] = \sum_{y_j \in R_Y} E[X|Y=y_j] P_Y(y_j) \quad \begin{matrix} P_X(x) = \sum_{y_j \in R_Y} P_{X|Y}(x|y_j) P_Y(y_j) \\ \downarrow \\ E[X] \end{matrix} \quad \begin{matrix} \sum_{y_j \in R_Y} E[X|Y=y_j] \\ \downarrow \\ E[X|Y=y_j] \end{matrix} \quad \begin{matrix} \{\text{no change.}\} \\ \downarrow \end{matrix}$$

### 5. Law of Total Variance

$$\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}(E[X|Y])$$



### 6. Law of the Unconscious Statistician.

$$E[g(X,Y)] = \sum_{y_j \in R_Y} \sum_{x_i \in R_X} g(x_i, y_j) P_{XY}(x_i, y_j) \quad \begin{matrix} E[g(x)] = \sum_{x_i \in R_X} g(x_i) P(X=x_i) \\ \downarrow \\ g(x_i) \end{matrix}$$

$$g(x_i, y_j) \quad P(X=x_i, Y=y_j)$$

### 7. Law of Iterated Expectation.

$$E[X] = E[E[X|Y]] \quad \leftarrow \text{This is kinda useful when you can find } f_{X|Y} \text{ but not } f_X$$

CHANGE PMF to PDF, and SUMMATION to INTEGRAL  
you will get continuous version of all the laws.

$$\text{eg. } E[X|Y=y] = \int_{-\infty}^{\infty} z f_{X|Y}(z|y) dz$$

$$E[g(x,y)] = \iint_{-\infty}^{\infty} g(x,y) f_{XY}(x,y) dx dy$$

If  $X$  and  $Y$  are INDEPENDENT variables,

- $E[XY] = E(X)E(Y)$ ; the converse is not always true.
  - $E[g(X)h(Y)] = E[g(X)]E[h(Y)]$ .
  - $E[X|Y] = EX$ ;
  - $E[g(X)|Y] = E[g(X)]$ .
- $\star$  all these are true for uncorrelated variables as independent  $\rightarrow$  uncorrelated!

$E[XY] = E[X]E[Y]$  does not imply independence.

(BUT IT DOES IMPLY UNCORRELATED!)

Let  $X_1, X_2, \dots, X_i$  be independent random variables and also independent of  $N$ . Let

$$Y = \sum_{i=1}^N X_i.$$

We have

$$\begin{aligned} E[Y] &= E[X]E[N] \\ Var(Y) &= E[N]Var(X) + (E[X])^2Var(N) \end{aligned}$$

} just memorise,  
quite hard.

## 6.6 Sum of Two Random Variables.

### 6.1 Method of Convolution

If we have two distributions, e.g.  $X \sim N(\mu_x, \sigma_x^2)$  and  $Y \sim N(\mu_y, \sigma_y^2)$  and we want to find the distribution of  $X+Y$ , e.g.  $X+Y \sim N(\dots)$  (more specific case check 6.6.3) we need to first find the PDF of  $X+Y$ :  $\star$  remember, PDF of  $X+Y$  is NOT just  $f_x(x) + f_y(y)$ !

let  $z = X+Y$ ,

$w$  is just a dummy variable!

If don't know  $X$  and  $Y$  are independent or not:

$$f_z(z) = \int_{-\infty}^{\infty} f_{X,Y}(w, z-w) dw \quad \text{or} \quad f_z(z) = \int_{-\infty}^{\infty} f_{X,Y}(z-w, w) dw$$

If know  $X$  and  $Y$  are independent:

$$f_z(z) = \int_{-\infty}^{\infty} f_X(w) f_Y(z-w) dw \quad \text{or} \quad f_z(z) = \int_{-\infty}^{\infty} f_X(z-w) f_Y(w) dw$$

$\star$  given in the data sheets:

Assuming independence of the added distributions:

$$K = I + J \Rightarrow p_K(k) = \underbrace{\sum_{j=0}^k p_I(k-j)p_J(j)}_{\text{discrete cases}}; \quad Z = X + Y \Rightarrow f_Z(z) = \int_{-\infty}^{+\infty} f_X(z-y) f_Y(y) dy$$

eg.  $X \sim \text{Exp}(\lambda)$ , find the distribution of  $2X$  (assume independent)

Step 1. Find the PDF of  $2X$ :

$$Z = X + Y \Rightarrow f_Z(z) = \int_{-\infty}^{+\infty} f_X(z-y)f_Y(y)dy$$

let  $Z=2X$ ,

$$\begin{aligned} f_Z(z) &= \int_{-\infty}^{\infty} f_X(z-w)f_X(w)dw \\ &\rightarrow = \int_0^z \lambda e^{-\lambda(z-w)} \cdot \lambda e^{-\lambda w} dw \\ &= \lambda^2 \int_0^z e^{-\lambda z + \lambda w} \cdot e^{-\lambda w} dw \\ &= \lambda^2 \int_0^z e^{-\lambda z} dw \\ &= \lambda^2 [we^{-\lambda z}]_0^z = \lambda^2 z e^{-\lambda z} \end{aligned}$$

we know that if  $X \sim \text{Exp}(\lambda)$   
 $f_X(x) = \lambda e^{-\lambda x}$

3. Gamma Distribution.  
 $X \sim \text{Gamma}(d, \beta)$  ( $d$  and  $\beta$  are shape and scale parameters)  
 $f(x) = \begin{cases} \frac{x^{d-1} e^{-x/\beta}}{\beta^d \Gamma(d)}, & \text{for } x > 0 \\ 0, & \text{otherwise} \end{cases}$   
 $E(X) = d\beta$ ;  $\text{Var}(X) = d\beta^2$

$$\begin{aligned} \beta &= \frac{1}{\lambda}, d = 2: f(z) = \frac{z^{2-1} e^{-z/\lambda}}{\frac{1}{\lambda}^2 (2-1)!} = \lambda^2 z e^{-\lambda z} \\ \therefore 2X &\sim \text{Gamma}(2, \frac{1}{\lambda}) \end{aligned}$$

### 6.6.2 Moments.

$X$  and  $Y$  doesn't have to be independent for this to be true.

$$E[aX+bY+c] = aE[X] + bE[Y] + c$$

$$\text{Var}(aX+bY+c) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

if only  $X$  and  $Y$  are independent,  $\text{Cov}(X, Y) = 0$  and hence:

$$\text{Var}(aX+bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y)$$

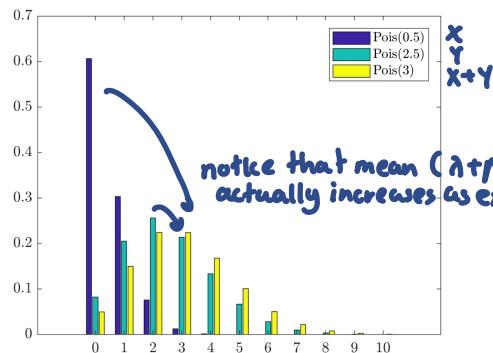
### 6.6.3 Distributions.

#### 1. Sum of Poisson Variables

If  $I \sim \text{Po}(\lambda)$ ,  $J \sim \text{Po}(\mu)$  are two independent discrete variables, then:



$$I+J \sim \text{Po}(\lambda+\mu)$$

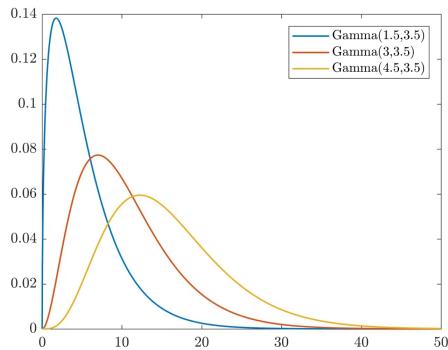


$$E(X) = \lambda + \mu$$

## 2. Sum of Gamma Variables

If  $X \sim \text{Gamma}(d_1, \beta)$  and  $Y \sim \text{Gamma}(d_2, \beta)$  are two independent variables, then:

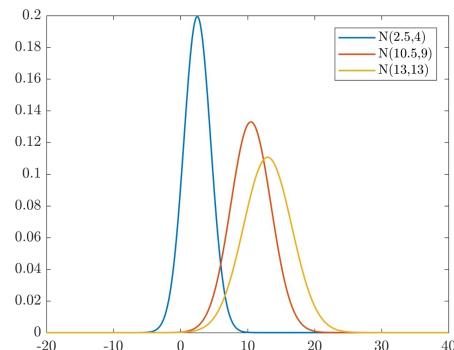
$$X + Y \sim \text{Gamma}(d_1 + d_2, \beta)$$



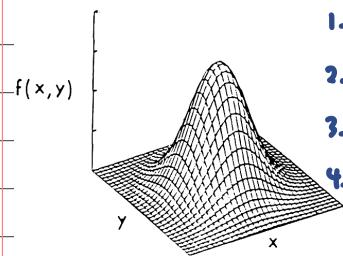
### 3. Sum of Normal Variables

If  $X \sim N(\mu_1, \sigma_1^2)$  and  $Y \sim N(\mu_2, \sigma_2^2)$  are two independent variables, then :

$$X+Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$



## 6.7 Bivariate Normal Distribution (Jointly Normal with 2 independent vars.)



- 1.** If  $X$  and  $Y$  are bivariate normal, then by letting  $a = 1$ ,  $b = 0$ , we conclude  $X$  must be normal.

**2.** If  $X$  and  $Y$  are bivariate normal, then by letting  $a = 0$ ,  $b = 1$ , we conclude  $Y$  must be normal.

**3.** If  $X \sim N(\mu_X, \sigma_X^2)$  and  $Y \sim N(\mu_Y, \sigma_Y^2)$  are independent, then they are jointly normal.

**4.** If  $X \sim N(\mu_X, \sigma_X^2)$  and  $Y \sim N(\mu_Y, \sigma_Y^2)$  are jointly normal, then

$$X + Y \sim N\left(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2 + 2\rho(X, Y)\sigma_X\sigma_Y\right)$$

properties of binormal  
**NOT IMPORTANT**.

1 and 2 basically just says  
X and Y are normal.

3 says X and Y are independent.

$$X + Y \sim N\left(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2 + 2\rho(X, Y)\sigma_X\sigma_Y\right).$$

This one is important! Memorise this instead of next page .

$$ax+by \sim N\left(E(ax+by), \text{Var}(ax+by)\right)$$

$$\sim N\left(aE(x)+bE(y), a^2\text{Var}(x)+b^2\text{Var}(y) + 2ab\text{Cov}(x,y)\right)$$

## 1. Bivariate Normal Distribution. ( $\mu_X, \mu_Y \in \mathbb{R}, \sigma_X, \sigma_Y > 0, \rho \in (-1, 1)$ and $\rho \neq 0$ )

**Definition:** Two random variables  $X$  and  $Y$  are said to have a **bivariate normal distribution** with parameters  $\mu_X, \sigma_X^2, \mu_Y, \sigma_Y^2$ , and  $\rho$ , if their joint PDF is given by

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x-\mu_X}{\sigma_X}\right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2 - 2\rho\frac{(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y}\right]\right\}$$

where  $\mu_X, \mu_Y \in \mathbb{R}, \sigma_X, \sigma_Y > 0$  and  $\rho \in (-1, 1)$  are all constants.

$$\left\{ \begin{array}{l} \text{set } \mu_X, \mu_Y = 0 \\ \text{and } \sigma_X, \sigma_Y = 1 \end{array} \right\} X \sim N(0, 1), Y \sim N(0, 1)$$

## 2. Standard Bivariate Normal Distribution with Correlation Coefficient

Two random variables  $X$  and  $Y$  are said to have the **standard bivariate normal distribution with correlation coefficient  $\rho$**  if their joint PDF is given by

$$f_{XY}(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}[x^2 - 2\rho xy + y^2]\right\}$$

where  $\rho \in (-1, 1)$ .

$$\left\{ \text{set } \rho = 0 \right.$$

## 3. Standard Bivariate Normal Distribution

If  $\rho = 0$ , then we just say that  $X$  and  $Y$  have the standard bivariate normal distribution.

$$f_{XY}(x, y) = \frac{1}{2\pi} \exp\left\{-\frac{1}{2}(x^2 + y^2)\right\}$$

Don't  
memorise!  
not  
important.

### Useful theorems:

1. Let  $X$  and  $Y$  be two bivariate normal random variables. Then there exist independent standard normal random variables  $Z_1$  and  $Z_2$  such that

$$\begin{cases} X = \sigma_X Z_1 + \mu_X \\ Y = \sigma_Y (\rho Z_1 + \sqrt{1 - \rho^2} Z_2) + \mu_Y \end{cases}$$

2. Suppose  $X$  and  $Y$  are jointly normal random variables with parameters  $\mu_X, \sigma_X^2, \mu_Y, \sigma_Y^2$ , and  $\rho$ . Then, given  $X = x$ ,  $Y$  is normally distributed with

$$E[Y|X = x] = \mu_Y + \rho \sigma_Y \frac{x - \mu_X}{\sigma_X},$$
$$Var(Y|X = x) = (1 - \rho^2) \sigma_Y^2.$$

3. If  $X$  and  $Y$  are bivariate normal and uncorrelated, then they are independent.

not important  
can skip

\* we always say if independent then it is correlated and NOT THE OTHER WAY but...

THIS IS A SPECIAL CASE.  
only bivariate normal and uncorrelated does mean independence !

## C7 Confidence Interval

(No notes on t-distribution and  $\chi^2$ -distribution cause we just need to know when to use it and how, and not the details of the distributions)

### 7.0 Pre-requisites.

Before we start with this chapter, let's be clear about samples and population.

Population is the entire data that you are interested / studying.

Samples is a subset of population that you are interested / studying.

In simple term, population is the "entire cake" and sample is just a few slices of it. (The slices of cakes, a.k.a the sample, can be used to make inference about the whole cake, a.k.a the population)



\* This chapter, is mainly about relating samples to population... which is very USEFUL! (we can't always know about the population, eg: you can't measure the average height of human population in the entire world, but you can measure some of it and make an inference about the entire population!)

### 7.1 Introduction to $E(\bar{x})$ and $\text{Var}(\bar{x})$

$X_1, X_2, X_3, \dots, X_n$  are INDIVIDUAL OBSERVATIONS!  
eg. height of one human from our sample.  
and each height, are independent from one another  
i.e. my height does not depend on yours.

for expectation,  $E(\bar{x})$ :

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

$E(X_i) = E(X)$  since they are taken from population

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

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

\* WHAT IS THE MEANING OF THIS?

This says that the mean, of mean of a sample is equal to the mean of population!

←  $E(\bar{x})$   
←  $E(X)$   
this 'n' is kind of confusing. It's number of individual observation from ONE sample and not the number of sample.

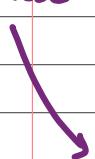
same for variance:

$$\text{Var}(\bar{x}) = \frac{\text{Var}(X_1 + X_2 + \dots + X_n)}{n^2} = \frac{\text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)}{n^2} = \frac{n\text{Var}(X)}{n^2}, \frac{\text{Var}(X)}{n}$$

$$\text{Var}(\bar{x}) = \frac{\text{Var}(X)}{n}$$

← This says that the variance, of the mean of a sample is equal to the variance of the population divided by n.

remember these



## 7.2 Limit Theorems.

as the name suggest, 'limit' theorems are theorem when  $\lim_{n \rightarrow \infty}$ , where sample size approaches infinity, which is the population!

### 1. Law of Large Numbers.

not really important  
BUT it is beneficial to understand

**Theorem:** Weak Law of Large Numbers

Let  $X_1, X_2, \dots, X_n$  be iid random variables with a finite expected value  $E(X_i) = \mu$ . Then, for any  $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|\bar{X} - \mu| \geq \epsilon) = 0$$

" $\bar{X} \rightarrow \mu$  if  $n \rightarrow \infty$ "

"limit as  $n$  approaches  $\infty$  of probability of sample mean minus population mean is greater than any positive numbers is 0"

in simple terms,  
the difference, between sample mean, and population mean is going to be equal to zero,  
WHEN the sample size,  $n$  approaches  $\infty$

### 2. Central Limit Theorem. (distribution of sample mean is ALWAYS NORMAL if sample size is big)

**Definition:**

Let  $X_1, X_2, \dots, X_n$  be iid random variables with expected value  $E(X_i) = \mu < \infty$  and variance  $0 < Var(X_i) = \sigma^2 < \infty$ . Then, the random variable

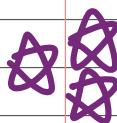
$$Z_n = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{X_1 + X_2 + \dots + X_n - n\mu}{\sqrt{n}\sigma}$$

converges to the standard normal random variable:

$$\lim_{n \rightarrow +\infty} P(Z_n \leq x) = \Phi(x), \quad \text{for all } x \in R$$

where  $\Phi(x)$  is the standard normal CDF.

basically it just means:



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

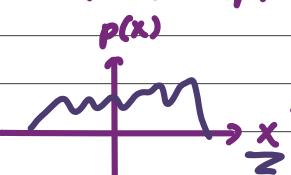
VERY IMPORTANT

after this, we can just treat it as a simple normal distribution question from c4.

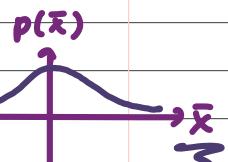
like we can standardize it to make finding probability easier  
where :

$$\begin{aligned} Z &= \frac{X - E(X)}{\sqrt{Var(X)}} & X \rightarrow \bar{X}, E(X) \rightarrow \mu \\ & & Var(X) \rightarrow \sigma^2/n \\ &= \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \end{aligned}$$

what is  $\bar{X} \sim N(\mu, \sigma^2)$  or sample distribution?



this is a distribution of a population.  
It can be anything.



if we take sample from the population and find the sample mean, and repeat, plot it on a graph, if  $n$  is big, we will get a normal distribution!

DOESN'T MATTER WHAT DISTRIBUTION THE POPULATION IS!

### 7.3 Confidence Interval

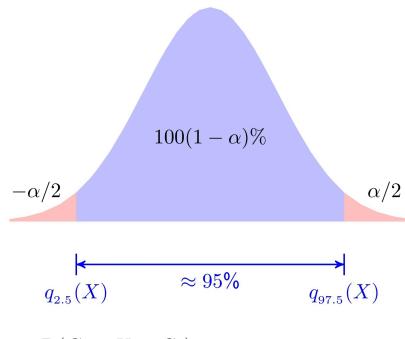
Confidence Interval of  $p\%$ . means you are  $p\%$ . confident that the outcome will be in your interval.

e.g. if we say 95% confidence interval of random variable  $X$  is between 5 and 12, it means 95% of the time,  $X$  will be in between 5 and 12.

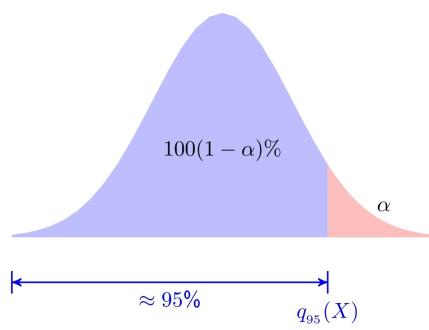
Although we can use confidence interval for a lot of case like (random variable  $X$ ), we will usually find CI of sample mean,  $\bar{X}$ ; and CI of sample variance,  $\hat{s}^2$ .  
The reason being in real life we usually don't know the distribution of the random variable  $X$  but we can estimate the distribution of  $\bar{X}$ , or  $\hat{s}^2$  ( $\sim N$ ,  $\sim t$ ,  $\sim \chi^2$  and etc.)

There are two types of CI, one-sided CI and two-sided CI:

Two-sided interval:



One-sided interval:



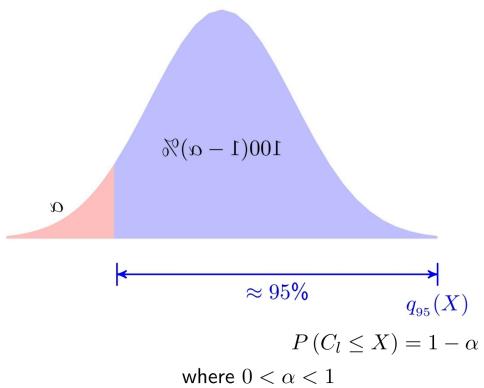
$$P(C_l \leq X \leq C_u) = 1 - \alpha \quad 0 < \alpha < 1$$

where  $C_l$  and  $C_u$  are the lower and upper confidence limits

$$P(X \leq C_u) = 1 - \alpha$$

where  $0 < \alpha < 1$

One-sided interval:



$$P(C_l \leq X) = 1 - \alpha$$

where  $0 < \alpha < 1$

e.g.

Confidence interval for a random variable.

Consider that the concentration of a substance  $X$  [mg/l] in a river is distributed as  $N(3, 0.25)$ . What is the value of the random variable  $X$  at a level of confidence 95%?

$$X \sim N(3, 0.25)$$

$$P(C_L \leq X \leq C_H) = 0.95$$

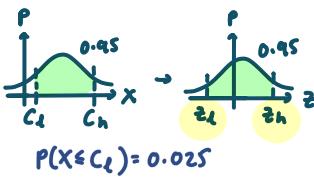
$$P(X \leq C_L) = 0.975$$

$$P\left(Z \leq \frac{C_H - 3}{\sqrt{0.25}}\right) = 0.975$$

$$\frac{C_H - 3}{\sqrt{0.25}} = 1.96$$

$$C_H = 1.96 \cdot \sqrt{0.25} + 3$$

$$C_H = 3.98$$



$$P(X \leq C_L) = 0.025$$

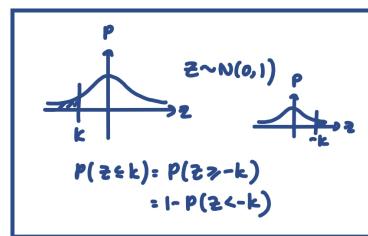
$$P\left(Z \leq \frac{C_L - 3}{\sqrt{0.25}}\right) = 0.025$$

$$\frac{C_L - 3}{\sqrt{0.25}} = -1.96$$

$$C_L = -1.96 \cdot \sqrt{0.25} + 3 \\ = 2.02$$

$$\therefore X \in (2.02, 3.98)$$

if doesn't mention anything  
→ two-sided CI



## 7.4 Sampling Distribution.

$$\bar{x}, s^2, \dots$$

Sampling Distribution → distribution of sample statistic

Sampling Distribution can be used to estimate population statistic

$$\mu, \sigma^2$$



We know that from Central Limit Theorem,  $\bar{X} \sim N(\mu, \sigma^2/n)$ , if size of each sample,  $n \rightarrow \infty$ . This is a type of sampling distribution! Distribution of Sample Mean!  
However,  $\bar{X} \sim N(\mu, \sigma^2/n)$  is only when  $n \rightarrow \infty$ , what if  $n$  is small? what if  $\sigma^2$  is not known? what distribution will we use?

### 7.4.1 Sampling Distribution of Sample Mean.

Sample Mean

basically we will try to use  $\sim N$   
If  $\sigma^2$  is known, if not known, try to estimate it  
ONLY IF not possible to estimate ( $n$  small)  
we will use t-test.

Standard Deviation (Population) is known

→ if was told  $X$  is normally distributed,  
no matter  $n$  is big ( $> 30$ ) or small,  
 $\bar{X} \sim N(\mu, \sigma^2/n)$

→ if  $n$  is big ( $> 30$ ),  $\bar{X} \sim N(\mu, \sigma^2/n)$

Standard Deviation (Population) is not known.

→ if  $n$  is small ( $n < 30$ ) or was told not to estimate  $\sigma$   
$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim t(n-1)$$
 WHAT IS T-DISTRIBUTION?  
CHECK EXAMPLE 3 BELOW!

→ if  $n$  is big ( $n > 30$ ), and is possible to estimate  $\sigma$   
(given dist. of  $\bar{X}$ ),  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$  with estimated  $\sigma$

\* for year 2 – statistics, don't worry about:

- s.d. known,  $X$  is not normally dist. AND  $n$  is small. (technically we use t-test)

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

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

## (two-sided, population s.d. known)

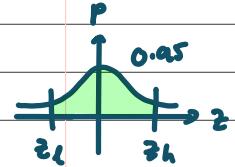
**eg.** Consider the compressive strengths of 40 test cubes of concrete. They have sample mean and standard deviation 60.14 and 5.02 N/mm<sup>2</sup>. Assuming a normal distribution and that the population standard deviation is well estimated by the sample standard deviation, estimate the 95% confidence level for the population mean.

$$\sigma^2 \text{ known, } \sigma^2 = s^2 = 5.02$$

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n}) \rightarrow \bar{x} \sim N(\mu, \frac{5.02}{40}) \rightarrow z = \frac{\bar{x} - \mu}{\sqrt{\sigma^2/n}} \sim N(0, 1)$$

$\mu$  is unknown, that's why we need to find the CI of  $\mu$ ! we want to know for 95% confidence, what range (interval) will  $\mu$  be at.

$$P(-z_{\alpha/2} < \frac{\bar{x} - \mu}{\sqrt{\sigma^2/n}} < z_{\alpha/2}) = 0.95$$



$$\text{from table, } P(-1.96 < z < 1.96) = 0.95:$$

$$P(-1.96 < \frac{\bar{x} - \mu}{\sqrt{\sigma^2/n}} < 1.96) = 0.95$$

$$P(\bar{x} - 1.96\sqrt{\sigma^2/n} < \mu < \bar{x} + 1.96\sqrt{\sigma^2/n}) = 0.95$$

$$\text{substituting } \bar{x} = 60.14, \sigma^2 = 5.02 \text{ and } n = 40,$$

$$P(58.59 < \mu < 61.71) = 0.95$$

$$\therefore \text{CI of 95% of } \mu = (58.59, 61.71)$$

## (one-sided, population s.d. unknown but can be estimated)

**eg. 2**

Consider the measurement of time-intervals  $x_i (i = 1, 2, \dots, 204)$  between vehicles at an observation point. The sample mean was found to be  $\bar{x} = 0.551$  mins and the exponential distribution provides a good fit to the data. The transport experts' concern is for small mean inter-arrival times, since they lead to congestion. Obtain a one-sided 99% confidence interval for the population mean.

→ s.d. (population) unknown, but able to estimate  $\sigma^2$

cause given  $X \sim \text{Exp}(\lambda)$ , so treat this question as s.d. known.

→  $n$  big, →  $\bar{x} \sim N(\mu, \sigma^2/n)$

how to estimate  $\sigma^2$ ?

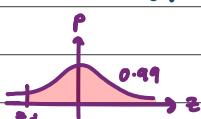
if  $X \sim \text{Exp}(\lambda)$

$$\left. \begin{array}{l} E(X) = 1/\lambda \\ \text{Var}(X) = 1/\lambda \end{array} \right\} E(X) = \text{Var}(X)$$

we know  $\bar{x} = 0.551$

$$E(X) \approx \bar{x} = 0.551 = 1/\lambda \quad \text{Var}(X) = \sigma^2 = \frac{1}{0.551}$$

$$\lambda = 1/0.551$$



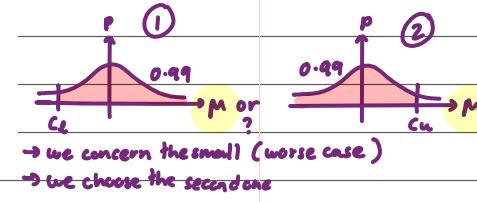
$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$P(z > z_d) = 0.99$$

$$P\left(\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} > z_d\right) = 0.99$$

$$P\left(\mu < \bar{x} - z_d \frac{\sigma}{\sqrt{n}}\right) = 0.99$$

$$P(\mu < 0.558) = 0.99$$

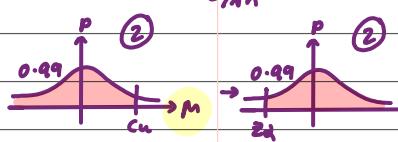


→ we concern the small (worse case)

→ we choose the second one

BUT! when convert to z...

$$\text{since } z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}, \text{ hence ...}$$



It will flip!  
hence we find  
 $P(z > z_d) = 0.99$ !

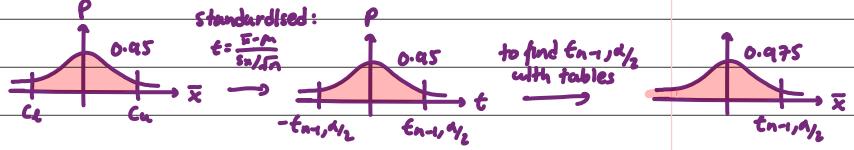
## two-sided, population s.d. unknown

Consider the compressive strengths of 40 test cubes of concrete. They have sample mean and standard deviation 60.14 and 5.02 N/mm<sup>2</sup>. Assuming a normal distribution and that the population standard deviation is well estimated by the sample standard deviation, estimate the 95% confidence level for the population mean.

Reconsider Example 43 with the compressive strengths of 40 test cubes of concrete.

Without assuming that  $\sigma$  is known estimate the 95% confidence level for the population mean.  
→ use t-test!

$$T = \frac{\bar{x} - \mu}{\hat{s}_x / \sqrt{n}} \sim t(n-1)$$



$$P(-t_{n-1, \alpha/2} < T < t_{n-1, \alpha/2}) = 0.95$$

$$P(-t_{n-1, \alpha/2} < \frac{\bar{x} - \mu}{\hat{s}_x / \sqrt{n}} < t_{n-1, \alpha/2}) = 0.95$$

$$P(\bar{x} - t_{n-1, \alpha/2} \cdot \frac{\hat{s}_x}{\sqrt{n}} < \mu < \bar{x} + t_{n-1, \alpha/2} \cdot \frac{\hat{s}_x}{\sqrt{n}}) = 0.95$$

$$P(60.14 - 2.023 \cdot \frac{5.02}{\sqrt{40}} < \mu < 60.14 + 2.023 \cdot \frac{5.02}{\sqrt{40}}) = 0.95$$

$$P(58.53 < \mu < 61.76) = 0.95$$

$$\therefore C.I. \text{ of } 95\% = (58.53, 61.76)$$

How to get t-values from t-distribution table?

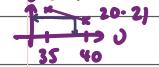
STUDENT'S t Table

ν 60.0 66.7% 75.0% 80.0% 87.5% 90.0% 95.0% 97.5% 99.0% 99.5% 99.9%

|    |       |       |       |       |       |       |       |        |        |        |        |
|----|-------|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|
| 1  | 0.325 | 0.577 | 1.000 | 1.376 | 2.414 | 3.078 | 6.314 | 12.706 | 31.821 | 63.657 | 318.31 |
| 2  | 0.289 | 0.500 | 0.691 | 1.061 | 1.604 | 1.886 | 2.920 | 4.303  | 6.965  | 9.925  | 23.237 |
| 3  | 0.277 | 0.476 | 0.765 | 0.978 | 1.423 | 1.638 | 2.353 | 3.182  | 4.541  | 5.841  | 10.215 |
| 4  | 0.271 | 0.464 | 0.741 | 0.941 | 1.384 | 1.582 | 2.132 | 2.776  | 3.747  | 4.604  | 7.173  |
| 5  | 0.265 | 0.452 | 0.720 | 0.904 | 1.329 | 1.520 | 2.070 | 2.700  | 3.482  | 4.302  | 6.393  |
| 6  | 0.260 | 0.453 | 0.718 | 0.896 | 1.273 | 1.440 | 1.943 | 2.447  | 3.143  | 3.700  | 5.208  |
| 7  | 0.263 | 0.449 | 0.711 | 0.896 | 1.254 | 1.415 | 1.895 | 2.365  | 2.994  | 3.499  | 4.785  |
| 8  | 0.262 | 0.447 | 0.706 | 0.889 | 1.240 | 1.397 | 1.860 | 2.306  | 2.895  | 3.355  | 4.501  |
| 9  | 0.261 | 0.445 | 0.703 | 0.883 | 1.230 | 1.383 | 1.833 | 2.262  | 2.821  | 3.250  | 4.297  |
| 10 | 0.260 | 0.444 | 0.700 | 0.879 | 1.221 | 1.372 | 1.812 | 2.228  | 2.764  | 3.169  | 4.144  |
| 11 | 0.260 | 0.443 | 0.697 | 0.876 | 1.214 | 1.363 | 1.796 | 2.201  | 2.718  | 3.106  | 4.025  |
| 12 | 0.259 | 0.442 | 0.693 | 0.873 | 1.209 | 1.356 | 1.782 | 2.179  | 2.681  | 3.055  | 3.930  |
| 13 | 0.259 | 0.441 | 0.694 | 0.870 | 1.204 | 1.350 | 1.771 | 2.160  | 2.650  | 3.012  | 3.852  |
| 14 | 0.258 | 0.440 | 0.692 | 0.868 | 1.200 | 1.345 | 1.761 | 2.145  | 2.624  | 2.977  | 3.787  |
| 15 | 0.258 | 0.439 | 0.691 | 0.866 | 1.197 | 1.347 | 1.753 | 2.131  | 2.602  | 2.947  | 3.733  |
| 16 | 0.258 | 0.439 | 0.689 | 0.865 | 1.194 | 1.337 | 1.746 | 2.120  | 2.582  | 2.921  | 3.686  |
| 17 | 0.257 | 0.438 | 0.688 | 0.864 | 1.191 | 1.334 | 1.740 | 2.110  | 2.562  | 2.895  | 3.646  |
| 18 | 0.257 | 0.438 | 0.688 | 0.863 | 1.189 | 1.330 | 1.736 | 2.101  | 2.542  | 2.868  | 3.619  |
| 19 | 0.257 | 0.438 | 0.688 | 0.861 | 1.187 | 1.328 | 1.729 | 2.093  | 2.539  | 2.861  | 3.579  |
| 20 | 0.257 | 0.437 | 0.687 | 0.860 | 1.185 | 1.325 | 1.725 | 2.086  | 2.528  | 2.845  | 3.552  |
| 21 | 0.257 | 0.437 | 0.686 | 0.859 | 1.183 | 1.323 | 1.721 | 2.080  | 2.518  | 2.831  | 3.527  |
| 22 | 0.256 | 0.437 | 0.686 | 0.858 | 1.182 | 1.321 | 1.717 | 2.074  | 2.508  | 2.819  | 3.505  |
| 23 | 0.256 | 0.436 | 0.685 | 0.858 | 1.180 | 1.319 | 1.714 | 2.069  | 2.500  | 2.807  | 3.485  |
| 24 | 0.256 | 0.436 | 0.685 | 0.857 | 1.179 | 1.318 | 1.711 | 2.064  | 2.492  | 2.797  | 3.467  |
| 25 | 0.256 | 0.436 | 0.684 | 0.856 | 1.178 | 1.316 | 1.708 | 2.060  | 2.485  | 2.787  | 3.450  |
| 26 | 0.256 | 0.435 | 0.684 | 0.856 | 1.177 | 1.315 | 1.706 | 2.056  | 2.479  | 2.779  | 3.435  |
| 27 | 0.256 | 0.435 | 0.684 | 0.855 | 1.176 | 1.314 | 1.703 | 2.052  | 2.473  | 2.771  | 3.421  |
| 28 | 0.256 | 0.435 | 0.684 | 0.855 | 1.175 | 1.313 | 1.701 | 2.048  | 2.467  | 2.763  | 3.408  |
| 29 | 0.256 | 0.435 | 0.684 | 0.854 | 1.174 | 1.312 | 1.699 | 2.045  | 2.462  | 2.756  | 3.396  |
| 30 | 0.256 | 0.435 | 0.683 | 0.853 | 1.173 | 1.311 | 1.697 | 2.040  | 2.459  | 2.749  | 3.385  |
| 31 | 0.256 | 0.434 | 0.683 | 0.852 | 1.171 | 1.306 | 1.695 | 2.039  | 2.458  | 2.732  | 3.349  |
| 32 | 0.255 | 0.434 | 0.681 | 0.851 | 1.167 | 1.303 | 1.684 | 2.024  | 2.432  | 2.704  | 3.307  |
| 33 | 0.255 | 0.434 | 0.680 | 0.850 | 1.165 | 1.301 | 1.679 | 2.014  | 2.412  | 2.690  | 3.281  |
| 34 | 0.255 | 0.433 | 0.679 | 0.849 | 1.164 | 1.299 | 1.676 | 2.009  | 2.403  | 2.678  | 3.261  |
| 35 | 0.255 | 0.433 | 0.673 | 0.848 | 1.163 | 1.297 | 1.673 | 2.004  | 2.396  | 2.668  | 3.245  |
| 36 | 0.254 | 0.433 | 0.674 | 0.848 | 1.162 | 1.296 | 1.671 | 2.000  | 2.390  | 2.660  | 3.232  |
| 37 | 0.253 | 0.433 | 0.674 | 0.842 | 1.156 | 1.282 | 1.665 | 1.993  | 2.380  | 2.656  | 3.090  |

interpolate for  $v=39$ :

$$20.80 - \frac{34-35}{40-35} \times (20.80 - 20.21) = 20.228 \approx 20.23$$



## 7.4.2 Sampling Distribution of Sample Variance

unlike sample mean, sample variance is more straightforward :

$$V = \frac{(n-1) \hat{s}_x^2}{s^2} \sim \chi^2(n-1) , \text{ no matter } n \text{ big or small...}$$

| 7.4.1 Sampling Distribution of Sample Mean:                                                                                               |                                                                                                                                                                          |
|-------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Sample Mean                                                                                                                               | basically we will try to use $t$ -test<br>if $\sigma$ is known, if not known, try to estimate $\sigma$ (n small)<br>but it's not possible to estimate $\sigma$ (n small) |
| Standard Deviation (Population) is known                                                                                                  |                                                                                                                                                                          |
| → if we told $X$ is normally distributed,<br>no matter $n$ is big ( $> 30$ ) or small,<br>$\bar{X} \sim N(\mu, \sigma^2/n)$               | → if $n$ is small ( $n < 30$ ) or not told to estimate $\sigma$<br>$\bar{X} \sim t(n-1)$ DISTRIBUTION?<br>$\bar{X} \sim N(\mu, \sigma^2/n)$                              |
| → if $n$ is big ( $> 30$ ), $\bar{X} \sim (\mu, \sigma^2/n)$                                                                              | → if $n$ is big ( $n > 30$ ), and is possible to estimate $\sigma$<br>(mean of $X$ ), $\bar{X} \sim N(\mu, \sigma^2/n)$ with estimated $\sigma$                          |
| for year 2 – statistics, don't worry about:<br>– $s$ d. known, $X$ is not normally dist. AND $n$ is small. (technically we use $t$ -test) |                                                                                                                                                                          |

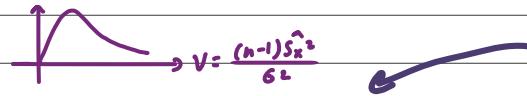
compared to sample mean...

eg.

Reconsider Example 39 with the compressive strengths of 40 test cubes of concrete.

Estimate the 99% confidence level for the population variance.

given that (in example 39),  $\bar{x} = 60.14$ ;  $\hat{s}_x = 5.02$



notice that here we don't write

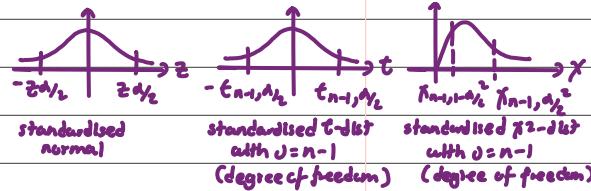
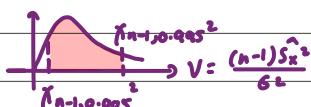
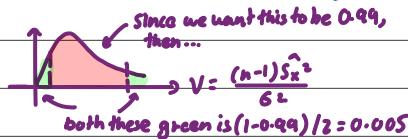
$$P(\chi_{n-1, d_{\alpha/2}}^2 < V < \chi_{n-1, 1-d_{\alpha/2}}^2)$$

cause  $\chi^2$  is not a symmetric dist. like N or t

$$P(\chi_{n-1, 1-d_{\alpha/2}}^2 < V < \chi_{n-1, d_{\alpha/2}}^2) = 0.99 \quad (\text{two-sided CI})$$

$$P\left(\frac{(n-1)\hat{s}_x^2}{6^2} < \frac{(n-1)\hat{s}_x^2}{6^2} < \chi_{n-1, d_{\alpha/2}}^2\right) = 0.99$$

$$P\left(\frac{n-1}{\chi_{n-1, d_{\alpha/2}}^2} \hat{s}_x^2 < \sigma^2 < \frac{n-1}{\chi_{n-1, 1-d_{\alpha/2}}^2} \hat{s}_x^2\right) = 0.99$$



\* I am taking  $v=40$  cause I am lazy to interpolate (It's fine)

to find  $\chi_{39, 0.005}^2$

CHI-SQUARED Table - 1/2

| $\nu$ | 0.1%   | 0.5%   | 1.0%   | 2.5%   | 5.0%   | 10.0%  | 12.5%  | 20.0%  | 25.0%  | 33.3%  | 50.0%  |
|-------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 1     | 0.000  | 0.000  | 0.001  | 0.004  | 0.016  | 0.025  | 0.064  | 0.102  | 0.186  | 0.455  |        |
| 2     | 0.002  | 0.010  | 0.020  | 0.051  | 0.103  | 0.211  | 0.466  | 0.575  | 0.811  | 1.386  |        |
| 3     | 0.024  | 0.072  | 0.115  | 0.216  | 0.352  | 0.584  | 0.692  | 1.005  | 1.213  | 1.568  | 2.366  |
| 4     | 0.091  | 0.207  | 0.297  | 0.484  | 0.711  | 1.064  | 1.219  | 1.649  | 1.923  | 2.378  | 3.357  |
| 5     | 0.210  | 0.412  | 0.554  | 0.831  | 1.145  | 1.610  | 1.808  | 2.343  | 2.675  | 3.216  | 4.351  |
| 6     | 0.381  | 0.676  | 0.872  | 1.237  | 1.635  | 2.204  | 2.441  | 3.070  | 3.455  | 4.074  | 5.348  |
| 7     | 0.598  | 0.989  | 1.239  | 1.690  | 2.167  | 2.833  | 3.106  | 3.822  | 4.255  | 4.945  | 6.346  |
| 8     | 0.857  | 1.344  | 1.646  | 2.180  | 2.733  | 3.490  | 3.797  | 4.594  | 5.071  | 5.826  | 7.344  |
| 9     | 1.152  | 1.735  | 2.088  | 2.700  | 3.325  | 4.168  | 4.507  | 5.380  | 5.898  | 6.716  | 8.343  |
| 10    | 1.479  | 2.156  | 2.558  | 3.247  | 3.940  | 4.865  | 5.234  | 6.179  | 6.737  | 7.612  | 9.342  |
| 11    | 1.834  | 2.603  | 3.053  | 3.816  | 4.574  | 5.578  | 5.975  | 6.989  | 7.584  | 8.514  | 10.341 |
| 12    | 2.214  | 3.074  | 3.571  | 4.404  | 5.226  | 6.304  | 6.729  | 7.807  | 8.436  | 9.420  | 11.340 |
| 13    | 2.617  | 3.565  | 4.107  | 5.009  | 5.892  | 7.042  | 7.493  | 8.634  | 9.299  | 10.331 | 12.340 |
| 14    | 3.041  | 4.075  | 4.660  | 5.629  | 6.571  | 7.790  | 8.266  | 9.467  | 10.165 | 11.245 | 13.339 |
| 15    | 3.483  | 4.601  | 5.229  | 6.262  | 7.261  | 8.547  | 9.048  | 10.307 | 11.037 | 12.163 | 14.339 |
| 16    | 3.942  | 5.142  | 5.812  | 6.908  | 7.962  | 9.312  | 9.812  | 11.512 | 11.912 | 13.083 | 15.338 |
| 17    | 4.416  | 5.697  | 6.408  | 7.564  | 8.672  | 10.085 | 10.633 | 12.002 | 12.792 | 14.006 | 16.338 |
| 18    | 4.905  | 6.265  | 7.015  | 8.231  | 9.390  | 10.865 | 11.435 | 12.857 | 13.675 | 14.931 | 17.338 |
| 19    | 5.407  | 6.844  | 7.633  | 8.797  | 10.117 | 11.651 | 12.242 | 13.716 | 14.562 | 15.859 | 18.338 |
| 20    | 5.921  | 7.434  | 8.260  | 9.591  | 10.851 | 12.443 | 13.055 | 14.578 | 15.452 | 16.788 | 19.337 |
| 21    | 6.447  | 8.034  | 8.897  | 10.283 | 11.591 | 13.240 | 13.873 | 15.445 | 16.344 | 17.720 | 20.337 |
| 22    | 6.983  | 8.643  | 9.542  | 10.982 | 12.338 | 13.441 | 14.695 | 16.314 | 18.653 | 21.337 |        |
| 23    | 7.529  | 9.260  | 10.196 | 11.689 | 13.091 | 14.845 | 15.521 | 17.187 | 18.137 | 19.587 | 22.337 |
| 24    | 8.085  | 9.886  | 10.856 | 12.401 | 13.849 | 15.659 | 16.351 | 18.062 | 19.037 | 20.523 | 23.337 |
| 25    | 8.649  | 10.520 | 11.524 | 13.120 | 14.611 | 16.473 | 17.184 | 18.940 | 19.939 | 21.461 | 24.337 |
| 26    | 9.222  | 11.160 | 12.198 | 13.844 | 15.379 | 17.292 | 18.021 | 19.820 | 20.843 | 22.399 | 25.336 |
| 27    | 9.803  | 11.808 | 12.879 | 14.573 | 16.151 | 18.114 | 18.861 | 20.703 | 21.749 | 23.339 |        |
| 28    | 10.391 | 12.461 | 13.565 | 15.308 | 16.928 | 18.939 | 19.704 | 21.588 | 22.657 | 24.280 | 27.336 |
| 29    | 10.986 | 13.121 | 14.256 | 16.047 | 17.708 | 19.768 | 20.550 | 22.475 | 23.567 | 25.222 | 28.336 |
| 30    | 11.588 | 13.787 | 14.953 | 16.718 | 18.493 | 20.599 | 21.399 | 23.364 | 24.478 | 26.165 | 29.335 |
| 31    | 14.688 | 17.192 | 18.509 | 20.569 | 22.465 | 24.797 | 25.678 | 27.836 | 29.054 | 30.894 | 34.336 |
| 32    | 7.9-6  | 20.707 | 22.164 | 24.433 | 26.509 | 29.051 | 30.008 | 32.345 | 33.660 | 35.643 | 39.335 |
| 33    | 21.251 | 24.611 | 25.901 | 28.366 | 30.612 | 33.350 | 34.379 | 36.884 | 38.291 | 40.407 | 44.335 |
| 34    | 24.674 | 27.991 | 29.707 | 32.357 | 34.764 | 37.689 | 38.175 | 41.449 | 42.942 | 45.184 | 49.355 |
| 35    | 28.173 | 31.735 | 33.570 | 36.398 | 38.954 | 42.069 | 43.220 | 46.406 | 47.610 | 49.972 | 54.335 |
| 36    | 31.738 | 35.534 | 37.485 | 40.482 | 43.188 | 46.459 | 47.680 | 50.641 | 52.294 | 54.770 | 59.335 |

CHI-SQUARED Table - 2/2

| $\nu$ | 60.0%  | 66.7%  | 75.0%  | 80.0%  | 87.5%  | 90.0%  | 95.0%  | 97.5%  | 99.0%  | 99.5%  | 99.9%  |
|-------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 1     | 0.708  | 0.936  | 1.323  | 1.642  | 2.354  | 2.706  | 3.841  | 5.024  | 6.635  | 7.879  | 10.828 |
| 2     | 1.833  | 2.197  | 2.773  | 3.219  | 4.159  | 4.605  | 5.991  | 7.378  | 9.210  | 10.597 | 13.816 |
| 3     | 2.946  | 3.405  | 4.108  | 4.642  | 5.739  | 6.251  | 7.815  | 9.348  | 11.345 | 12.838 | 16.266 |
| 4     | 4.045  | 4.579  | 5.385  | 5.989  | 7.214  | 7.779  | 9.488  | 11.143 | 13.277 | 14.860 | 18.467 |
| 5     | 5.132  | 5.730  | 6.626  | 7.289  | 8.625  | 9.236  | 11.070 | 12.833 | 15.086 | 16.750 | 20.515 |
| 6     | 6.211  | 6.867  | 7.841  | 8.558  | 9.992  | 10.645 | 12.592 | 14.449 | 16.812 | 18.548 | 22.458 |
| 7     | 7.283  | 7.999  | 9.037  | 9.808  | 11.326 | 12.017 | 14.067 | 16.013 | 18.475 | 20.278 | 24.322 |
| 8     | 8.351  | 9.107  | 10.219 | 11.030 | 12.630 | 13.369 | 15.507 | 17.535 | 20.090 | 21.955 | 26.125 |
| 9     | 9.414  | 10.215 | 11.389 | 12.242 | 13.926 | 14.684 | 16.919 | 19.023 | 21.666 | 23.589 | 27.877 |
| 10    | 10.473 | 11.317 | 12.549 | 13.442 | 15.193 | 15.987 | 18.307 | 20.483 | 23.209 | 25.188 | 29.588 |
| 11    | 11.530 | 12.414 | 13.701 | 14.631 | 16.457 | 17.275 | 19.675 | 21.920 | 24.725 | 26.757 | 31.264 |
| 12    | 12.584 | 13.506 | 14.845 | 15.812 | 17.703 | 18.549 | 21.026 | 23.337 | 26.217 | 28.300 | 32.910 |
| 13    | 13.636 | 14.593 | 15.984 | 16.935 | 18.919 | 19.812 | 22.362 | 24.736 | 27.688 | 29.819 | 34.528 |
| 14    | 14.685 | 15.680 | 17.117 | 18.151 | 20.166 | 21.064 | 23.685 | 26.211 | 29.119 | 31.191 | 36.123 |
| 15    | 15.733 | 16.761 | 18.245 | 19.311 | 21.384 | 22.307 | 24.996 | 27.488 | 30.578 | 32.801 | 37.697 |
| 16    | 16.780 | 17.840 | 19.369 | 20.465 | 22.595 | 23.542 | 26.296 | 28.845 | 32.000 | 34.267 | 39.252 |
| 17    | 17.824 | 18.917 | 20.489 | 21.615 | 23.799 | 24.769 | 27.587 | 30.191 | 33.409 | 35.718 | 40.790 |
| 18    | 18.868 | 19.991 | 21.605 | 22.760 | 24.997 | 25.988 | 28.869 | 31.526 | 34.805 | 37.156 | 42.312 |
| 19    | 19.910 | 21.063 | 22.718 | 23.900 | 26.189 | 27.204 | 30.244 | 32.852 | 36.191 | 38.582 | 43.820 |
| 20    | 20.951 | 22.133 | 23.828 | 25.038 | 27.376 | 28.412 | 31.410 | 34.170 | 37.566 | 39.997 | 45.315 |
| 21    | 21.991 | 23.201 | 24.935 | 26.171 | 28.559 | 29.615 | 32.671 | 35.479 | 38.932 | 41.401 | 46.797 |
| 22    | 23.031 | 24.268 | 26.039 | 27.301 | 29.737 | 30.813 | 33.924 | 36.781 | 40.289 | 42.796 | 48.268 |
| 23    | 24.069 | 25.333 | 27.141 | 28.429 | 30.911 | 32.007 | 35.172 | 38.076 | 41.638 | 44.181 | 49.728 |
| 24    | 25.106 | 26.397 | 28.241 | 29.553 | 32.081 | 33.191 | 36.415 | 39.364 | 42.980 | 45.559 | 51.179 |
| 25    | 26.143 | 27.459 | 29.339 | 30.675 | 33.247 | 34.382 | 37.652 | 40.646 | 44.314 | 46.928 | 52.620 |
| 26    | 27.179 | 28.520 | 30.435 | 31.795 | 34.410 | 35.563 | 38.885 | 41.923 | 45.642 | 48.290 | 54.052 |
| 27    | 28.214 | 29.580 | 31.528 | 32.912 | 35.570 | 36.741 | 40.113 | 43.195 | 46.963 | 49.645 | 55.476 |
| 28    | 29.249 | 30.639 | 32.620 | 34.037 | 36.276 | 37.916 | 41.337 | 44.466 | 48.278 | 50.993 | 56.892 |
| 29    | 30.283 | 31.697 | 33.711 | 35.139 | 37.881 | 39.487 | 42.557 | 45.722 | 49.588 | 52.336 | 58.301 |
| 30    | 31.311 | 32.754 | 34.800 | 36.250 | 39.033 | 40.256 | 43.773 | 46.979 | 50.892 | 53.672 | 59.703 |
| 31    | 32.367 | 34.724 | 36.800 | 38.250 | 41.033 | 42.659 | 46.099 | 49.802 | 53.203 | 57.342 | 66.619 |
| 32    | 33.422 | 35.804 | 38.024 | 40.223 | 41.778 | 44.755 | 49.059 | 49.802 | 53.203 | 57.342 | 6      |

# c8 Statistical Testing (Hypothesis Testing)

## 8.1 Introduction

We do statistical testing when we have a hypothesis in mind, and want to check if our hypothesis is true or not. (eg. mean of height of year 2's student = 170 cm)

To do hypothesis testing, we will implement the following steps (let  $\mu$  = mean of height of year 2's student)

1. null hypothesis,  $H_0 : \mu = 170$

(null hypothesis doesn't necessarily is our initial claim.)

In this case, our initial claim is " $\mu = 170$ ". What if my initial claim is "Year 2 average height is greater than Year 1"? Then my null hypothesis,  $H_0$  will be:

2. alternative hypothesis,  $H_1 :$   
 $\mu \neq 170$   
 (since we want to test if  $\mu$  is really 170, and not like whether it is greater or smaller than 170, it is a two-sided hypothesis test)

$\mu_2 = \mu_1$ , and alternative hypothesis will be  $\mu_2 \neq \mu_1$ .

This prove that initial claim is not always  $H_0$  or  $H_1$ .

Null hypothesis is just always the "neutral" one, and is always treated as true, unless proven false)

## 3. Define a statistical test

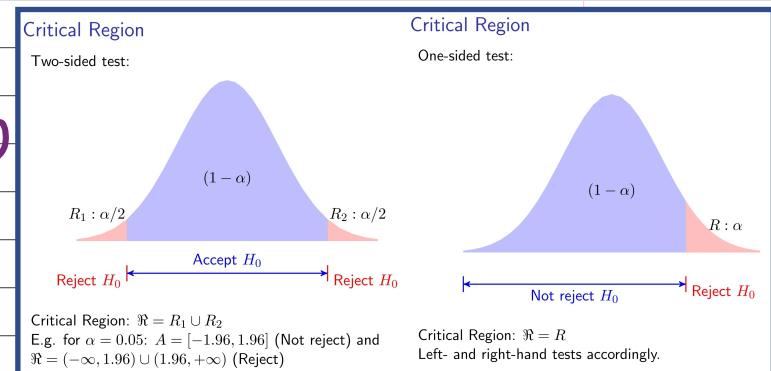
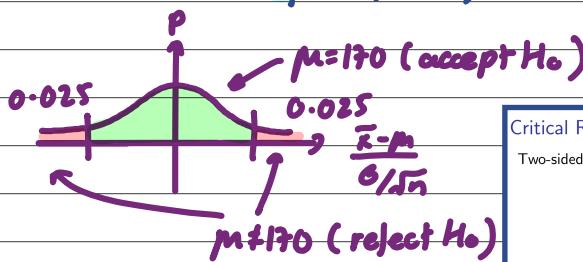
eg. If mean height of year 2's student is normal distribution (this is true if  $\sigma$  is known and  $n > 30$ ):

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

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

## 4. Define critical region OR find the probability of alternative hypothesis.

eg. if I want to test at 95% confidence (5% significant) whether  $\mu = 170$ , we need to check if  $P(\mu \neq 170) < 0.05$ , if yes, we reject  $H_0$  ( $H_1$  is true, i.e.  $\mu \neq 170$ )



eg. (just an example on how hypothesis testing procedure works)

Check if a coin is fair. Toss it 100 times and answer with 95% confidence. (5% significant)

1. Define  $H_0$ : Coin is fair (i.e.  $P(H) = 0.5$ )

probability of getting heads = 0.5 (mathematical term)

2. Define  $H_1$ : Coin is not fair (i.e.  $P(H) \neq 0.5$ )

3. Define a statistical test.

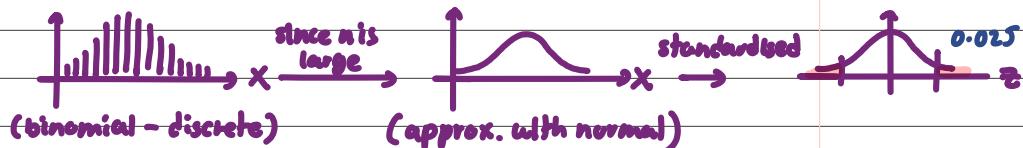
AT THIS STAGE (when defining a statistical test) WE MUST TREAT  $H_0$  AS TRUE!

i.e.  $P(H) = 0.5$ !

→ since tossing coins ( $n=100$ ) times, and each toss is ( $p=0.5$ ), this is BINOMIAL

→  $X \sim B(100, 0.5)$

→ to reject  $H_0$ ,



$$X \sim B(100, 0.5) \rightarrow X \sim N(50, 25)$$

$$E(X) = np = 50$$

$$\text{Var}(X) = np(1-p) = 25$$

$$\frac{X-50}{\sqrt{25}} \sim N(0,1)$$

we will reject  $H_0$  if  $P(z < -z_{\alpha/2}) = 0.025$

$$P(z > z_{\alpha/2}) = 0.025$$

$$P(z < -z_{\alpha/2}) = 0.025$$

$$P\left(\frac{X-50}{5} < -1.96\right) = 0.025 \quad P\left(\frac{X-50}{5} > 1.96\right) = 0.025$$

$$P(X < 40.2) = 0.025$$

$$P(X > 59.8) = 0.025$$



hence we will reject  $H_0$  (coin is not fair), if:  
we get 40 or less head out of 100 throws, or  
60 or more head out of 100 throws

we will accept  $H_0$  (coin is fair), if:  
we get 41 to 59 head out of 100 throws.

## 8.2 Asymmetric Case (One-sample Analysis) and Symmetric Case (Two-sample Analysis)

### Statistical Testing (Hypothesis Testing)

#### Asymmetric Case (one-sample)

Compare ONE SAMPLE's statistic to a known population statistic

this is a constant.

e.g. testing whether Year 2's average height is greater than world's average height

$$H_0: \mu_2 = 170 \quad (\mu_2 = \mu_{world} = 170)$$

$$H_1: \mu_2 > 170 \quad (\mu_2 > \mu_{world} = 170)$$

notice population mean is constant. That is why we call it one-sample analysis because only one sample, the other one is population and is constant.

Suppose A be the population with known mean:

#### 1. $\sigma_B$ is known ( $\sigma_B = \sigma_A = \sigma$ )

$$H_0: \mu_B = \mu_A \leftarrow \mu_A \text{ is constant}$$

$$H_1: \mu_B > \mu_A \text{ (one-sided)}$$

$$\mu_B \neq \mu_A \text{ (two-sided)}$$

$$z = \frac{\bar{x}_B - \mu_B}{\sigma_B / \sqrt{n}} \sim N(0, 1)$$

#### 2. $\sigma_B$ is not known (and can't be estimated)

$$H_0: \mu_B = \mu_A \leftarrow \mu_A \text{ is constant.}$$

$$H_1: \mu_B > \mu_A \text{ (one-sided)}$$

$$\mu_B \neq \mu_A \text{ (two-sided)}$$

$$T = \frac{\bar{x}_B - \mu_B}{\hat{s}_{x_B} / \sqrt{n}} \sim t(n-1)$$

#### Symmetric Case (two-sample)

Compare TWO SAMPLE's statistic

e.g. testing whether Year 2's average height is greater than Year 1's average height.

$$H_0: \mu_2 = \mu_1$$

$$H_1: \mu_2 > \mu_1$$

notice that both sample mean is NOT constant. That is why we called it two sample analysis.

#### 1. $\sigma_B$ is known, $\sigma_A$ is known

$$H_0: \mu_B = \mu_A \leftarrow \mu_B \text{ and } \mu_A \text{ both NOT constant}$$

$$H_1: \mu_B > \mu_A \text{ (one-sided)}$$

$$\mu_B \neq \mu_A \text{ (two-sided)}$$

$$z = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{\sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}} \sim N(0, 1)$$

#### 2. $\sigma_B$ and $\sigma_A$ is not known but ( $\sigma_A = \sigma_B$ )

$$H_0: \mu_B = \mu_A \leftarrow \mu_B \text{ and } \mu_A \text{ both NOT constant!}$$

$$H_1: \mu_B > \mu_A \text{ (one-sided)}$$

$$\mu_B \neq \mu_A \text{ (two-sided)}$$

$$t = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{\sqrt{\frac{\hat{s}_{x_A}^2}{n_A} + \frac{\hat{s}_{x_B}^2}{n_B}}} \sim t(n_A + n_B - 2)$$

$$\hat{s}_{x_B}^2 = \frac{(n_A - 1)\hat{s}_{x_A}^2 + (n_B - 1)\hat{s}_{x_B}^2}{n_A + n_B - 2}$$

pooled variance

#### 3. $\sigma_B$ and $\sigma_A$ is not known and ( $\sigma_A \neq \sigma_B$ )

same as (2) but  $\sim t(v)$  where  $v = \frac{\left(\frac{\hat{s}_{x_A}^2}{n_A} + \frac{\hat{s}_{x_B}^2}{n_B}\right)^2}{\frac{1}{n_A-1}\left(\frac{\hat{s}_{x_A}^2}{n_A}\right)^2 + \frac{1}{n_B-1}\left(\frac{\hat{s}_{x_B}^2}{n_B}\right)^2}$

### 8.3 Test Error

Two types of errors when testing:

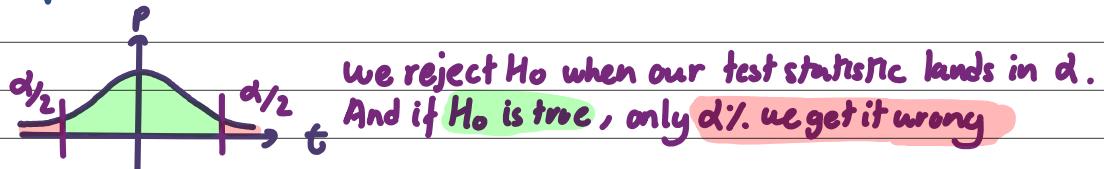
| Condition / Decision | accept $H_0$ | reject $H_0$ |
|----------------------|--------------|--------------|
| $H_0$ true           | $1 - \alpha$ | $\alpha$     |
| $H_1$ true           | $\beta$      | $1 - \beta$  |

Type 1

Type 2

Type 1 error is when  $H_0$  (Null Hypothesis) is TRUE, but we reject it.

The probability of this happening is just  $\alpha$  (significant level)



Type 2 error is when  $H_0$  (Null Hypothesis) is FALSE, but we accept it.

The probability of this happening is  $\beta$

$$\therefore \alpha = P(\text{reject } H_0 \mid H_0)$$

$$\therefore \beta = P(\text{accept } H_0 \mid H_1)$$

Two types of errors when testing:

| Condition / Decision | accept $H_0$ | reject $H_0$ |
|----------------------|--------------|--------------|
| $H_0$ true           | $1 - \alpha$ | $\alpha$     |
| $H_1$ true           | $\beta$      | $1 - \beta$  |

confidence level

significance

power

# VERY HARD EXAMPLE (one-sample, type 2 error, two-sample)

## 3. Samples of concrete beams

The densities of 30 concrete beams produced by a given supplier A have been tested.

In standard units, the sample mean is 20 and the sample standard deviation (using the unbiased formula) is 1.8.

A set of 10 beams is found in a warehouse, and the mean density of this new sample is 18 standard units.

c - Assuming that the sample mean and standard deviation found for the first sample of 30 beams can be used as reliable estimates of, respectively, the population mean and standard deviation of beams produced by supplier A, test with 5% significance the hypothesis that these 10 beams, although their mean density is smaller than 20, are

produced by supplier A.

$$\left\{ \begin{array}{l} n_A = 30, \bar{x}_A = 20, \hat{s}_{x_A} = 1.8 \end{array} \right.$$

$$\left\{ \begin{array}{l} n = 10, \bar{x} = 18 \end{array} \right.$$

$$\left\{ \begin{array}{l} \mu_A = \bar{x}_A = 20 \\ \sigma_A = \hat{s}_{x_A} = 1.8 \end{array} \right.$$

## Step 1. Define $H_0$ and $H_1$

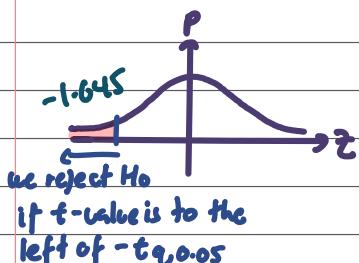
$$H_0: \mu = \mu_A (= 20)$$

$$H_1: \mu < \mu_A (= 20)$$

## Step 2. Define a statistical test $H_0$ IS TRUE (until proven false otherwise)

- using z-test as  $\sigma^2$  known ( $\sigma^2 = \sigma_A^2$ )  $\therefore \mu = \mu_A = 20$

$$\bar{z} = \frac{\bar{x} - \mu}{\hat{s}_{\bar{x}} / \sqrt{n}} \sim N(0,1)$$



to find what  $\bar{z}$  will give critical z-value ( $-z_{0.05}$ )

$$P(\bar{z} < -z_{0.05}) = 0.05$$

$$P(\bar{z} > z_{0.05}) = 0.05$$

$$P(\bar{z} < z_{0.05}) = 0.95$$

$$z_{0.05} = 1.645$$

$$\therefore P\left(\frac{\bar{x}-20}{1.8/\sqrt{10}} < -1.645\right) = 0.05$$

$$P(\bar{x} < 19.06) = 0.05$$

since  $\bar{x} = 18 < 19.06$ , reject  $H_0$ ! (these 10 beams are NOT produced by supplier A)

d - The suspicion arises that these 10 beams are actually produced by another supplier, B, for whom the population mean is estimated to be 19 standard units (we assume for this question that the population standard deviation is the same as for supplier A, i.e. 1.8). If the test carried out in (c) is now considered as a test of supplier A (null hypothesis) against supplier B (alternative hypothesis), what is the power of this test?

from part (c):

$$\mu_A = 20, \sigma_A = 1.8$$

from part (d):

$$\mu_B = 19, \sigma_B = 1.8$$

$$\begin{aligned} H_0: \mu = \mu_A (= 20) & \quad \text{power} = 1 - \beta = P(\text{reject } H_0 | H_1) \\ H_1: \mu = \mu_B (= 19) \end{aligned}$$

$$\text{power} = P(\text{reject } H_0 \mid H_1)$$

from part (c) we knew we reject  $H_0$  when  $\bar{z} < 19.0637$

$$\text{power} = P(\bar{z} < 19.0637 \mid \mu=19)$$

$$= P\left(z < \frac{19.0637 - 19}{1.8/\sqrt{10}}\right)$$

$$= P(z < 0.1119)$$

$$= 0.5445$$

e - It is now no longer assumed that the original sample of 30 beams provides a reliable estimate of the population statistics of beam densities from supplier A. Given that the sample standard deviation (unbiased) for the new sample of size 10, is 1.75 standard units, test with 5% significance whether this sample of 10 beams has a different mean from the original sample of 30 beams (in this test, the alternative hypothesis is, as in

(d), that the sample of 10 beams is from supplier B). Note that the two sample standard deviations are very similar.

$$\begin{cases} H_0: \mu_A = \mu_B \\ H_1: \mu_A \neq \mu_B \end{cases} \quad \begin{array}{l} \text{both } \mu_A, \mu_B \text{ are unknown (two-sample analysis)} \\ \sigma_A^2 \text{ and } \sigma_B^2 \text{ unknown but } (\sigma_A^2 = \sigma_B^2) \end{array}$$

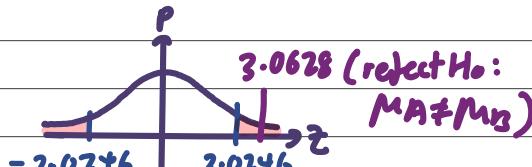
$$\mu_A - \mu_B = 0$$

$$t = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{\sqrt{\frac{s_p^2}{n_A} + \frac{s_p^2}{n_B}}} \sim t(n_A + n_B - 2)$$

$s_p^2 = \frac{(n_A-1)s_{x_A}^2 + (n_B-1)s_{x_B}^2}{n_A + n_B - 2}$   
pooled variance

$$s_p^2 = \frac{(30-1) \times 1.8^2 + (10-1) \times 1.75^2}{30+10-2} = 3.198$$

$$t = \frac{(20-18)-(0)}{\sqrt{\frac{3.198}{30} + \frac{3.198}{10}}} = 3.0628 \sim t(30+10-2)$$



$$P(T > t_{38, 0.025}) = 0.025$$

$$P(T < -t_{38, 0.025}) = 0.025$$

| $\nu$ | 60.0% | 66.7% | 75.0% | 80.0% | 87.5% | 90.0% | 95.0% | 97.5%  | 99.0%  | 99.5%  | 99.9%  |
|-------|-------|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|
| 1     | 0.325 | 0.577 | 1.000 | 1.376 | 2.414 | 3.078 | 6.311 | 12.706 | 31.821 | 63.657 | 318.31 |
| 2     | 0.289 | 0.500 | 0.816 | 1.061 | 1.604 | 1.886 | 2.929 | 4.303  | 6.965  | 9.925  | 22.327 |
| 3     | 0.277 | 0.476 | 0.767 | 0.978 | 1.423 | 1.638 | 2.353 | 3.182  | 4.541  | 5.841  | 10.215 |
| 4     | 0.271 | 0.464 | 0.741 | 0.941 | 1.344 | 1.533 | 2.132 | 2.776  | 3.747  | 4.604  | 7.173  |
| 5     | 0.267 | 0.457 | 0.727 | 0.920 | 1.301 | 1.476 | 2.015 | 2.571  | 3.365  | 4.032  | 5.893  |
| 6     | 0.265 | 0.453 | 0.718 | 0.906 | 1.273 | 1.440 | 1.945 | 2.447  | 3.143  | 3.707  | 5.208  |
| 7     | 0.263 | 0.449 | 0.711 | 0.896 | 1.254 | 1.415 | 1.895 | 2.365  | 2.998  | 3.499  | 4.785  |
| 8     | 0.262 | 0.447 | 0.703 | 0.889 | 1.240 | 1.397 | 1.860 | 2.305  | 2.896  | 3.355  | 4.501  |
| 9     | 0.261 | 0.445 | 0.701 | 0.883 | 1.230 | 1.383 | 1.833 | 2.262  | 2.821  | 3.250  | 4.297  |
| 10    | 0.260 | 0.444 | 0.700 | 0.879 | 1.221 | 1.372 | 1.812 | 2.228  | 2.764  | 3.169  | 4.144  |
| 11    | 0.260 | 0.443 | 0.697 | 0.876 | 1.214 | 1.363 | 1.796 | 2.209  | 2.718  | 3.106  | 4.025  |
| 12    | 0.259 | 0.442 | 0.695 | 0.873 | 1.209 | 1.356 | 1.782 | 2.179  | 2.681  | 3.055  | 3.930  |
| 13    | 0.259 | 0.441 | 0.699 | 0.870 | 1.204 | 1.350 | 1.771 | 2.165  | 2.650  | 3.012  | 3.852  |
| 14    | 0.258 | 0.440 | 0.692 | 0.868 | 1.200 | 1.345 | 1.761 | 2.145  | 2.624  | 2.977  | 3.787  |
| 15    | 0.258 | 0.439 | 0.689 | 0.865 | 1.196 | 1.341 | 1.752 | 2.125  | 2.595  | 2.949  | 3.733  |
| 16    | 0.258 | 0.439 | 0.689 | 0.860 | 1.194 | 1.338 | 1.743 | 2.109  | 2.583  | 2.921  | 3.696  |
| 17    | 0.257 | 0.438 | 0.689 | 0.863 | 1.191 | 1.333 | 1.740 | 2.110  | 2.567  | 2.898  | 3.646  |
| 18    | 0.257 | 0.438 | 0.688 | 0.862 | 1.189 | 1.330 | 1.734 | 2.101  | 2.552  | 2.878  | 3.610  |
| 19    | 0.257 | 0.438 | 0.688 | 0.861 | 1.187 | 1.328 | 1.729 | 2.093  | 2.539  | 2.861  | 3.579  |
| 20    | 0.257 | 0.437 | 0.687 | 0.860 | 1.185 | 1.325 | 1.725 | 2.085  | 2.528  | 2.845  | 3.552  |
| 21    | 0.257 | 0.437 | 0.686 | 0.859 | 1.182 | 1.322 | 1.721 | 2.080  | 2.518  | 2.831  | 3.537  |
| 22    | 0.256 | 0.437 | 0.686 | 0.858 | 1.182 | 1.321 | 1.717 | 2.074  | 2.508  | 2.819  | 3.505  |
| 23    | 0.256 | 0.436 | 0.685 | 0.858 | 1.180 | 1.319 | 1.714 | 2.069  | 2.500  | 2.807  | 3.485  |
| 24    | 0.256 | 0.436 | 0.685 | 0.857 | 1.179 | 1.318 | 1.711 | 2.064  | 2.492  | 2.797  | 3.467  |
| 25    | 0.256 | 0.436 | 0.684 | 0.856 | 1.178 | 1.316 | 1.708 | 2.060  | 2.485  | 2.787  | 3.450  |
| 26    | 0.256 | 0.436 | 0.684 | 0.856 | 1.177 | 1.315 | 1.706 | 2.059  | 2.479  | 2.779  | 3.435  |
| 27    | 0.256 | 0.435 | 0.684 | 0.855 | 1.176 | 1.314 | 1.703 | 2.053  | 2.473  | 2.771  | 3.421  |
| 28    | 0.256 | 0.435 | 0.684 | 0.854 | 1.175 | 1.313 | 1.701 | 2.048  | 2.467  | 2.763  | 3.408  |
| 29    | 0.256 | 0.435 | 0.683 | 0.854 | 1.174 | 1.311 | 1.699 | 2.045  | 2.462  | 2.756  | 3.396  |
| 30    | 0.256 | 0.435 | 0.683 | 0.854 | 1.173 | 1.310 | 1.697 | 2.041  | 2.457  | 2.750  | 3.385  |
| 35    | 0.255 | 0.434 | 0.682 | 0.852 | 1.170 | 1.306 | 1.696 | 2.030  | 2.438  | 2.724  | 3.340  |
| 40    | 0.255 | 0.434 | 0.681 | 0.851 | 1.167 | 1.303 | 1.685 | 2.021  | 2.423  | 2.704  | 3.307  |

$$2.030 - (2.030 - 2.021) \times \frac{38-35}{40-35}$$

$$= 2.0246$$

## 8.4 Goodness of Fit Test.

Goodness of fit test is to check how well a distribution fit onto data sample.  
In CS3 Goodness of Fit, we covered the following test

Tools we already have to assess the goodness of fit:

- ▶ qq-plots
- ▶ Compare CDFs
- ▶ Compare PDFs (use of histograms)

All these tools are visual ~ qualitative.

*not very good.*

two type of test covered in year 2:

Quantitative tests: *> better!*

- ▶ Decide with a level of confidence  $1 - \alpha$  if distribution  $D$  provides a good fit to the data

| Feature           | Chi-Squared                                      | Kolmogorov-Smirnov (KS)         |
|-------------------|--------------------------------------------------|---------------------------------|
| Data Type         | Binned Continuous                                | Continuous                      |
| Basis             | Observed vs. Expected Frequencies                | ECDF vs. Theoretical CDF        |
| Comparison Metric | Sum of squared relative differences ( $\chi^2$ ) | Maximum absolute difference (D) |

### 8.4.1 Chi-squared Test ( $\chi^2$ ) *& $\chi^2$ test is always a one-sided (specifically right-sided) test!*

1. Null Hypothesis  $H_0$ : Distribution is a good fit
2. Alternative Hypothesis  $H_1$ : Distribution is not a good fit
3. Test statistic  $C_n = \sum_{i=1}^l \frac{(O_i - E_i)^2}{E_i}$
4. Distribution of test statistic  $C_n \sim \chi^2(l - 1 - k)$
5. Critical Region  $\mathcal{R}$   $C_n > \chi^2_{l-1-k, 1-\alpha}$  (e.g.  $\alpha = 5\%$ )
6. Evaluate  $C_n$  under  $H_0$  if  $c_{n0} > \chi^2_{l-1-k, \alpha} \Rightarrow$  Reject  $H_0$
7.  $p$ -value  $p(C_n > c_{n0})$

*}  $H_0$  and  $H_1$  will be always like this!*

Expected frequencies:

$$E_i = (F_X(x_{i+1}) - F_X(x_i)) \times n \quad \forall i \in [2, l-1]$$

where  $F_X$  is the CDF under test.

For class 1:  $E_1 = F_X(x_2) \times n$

For class  $l$ :  $E_l = (1 - F_X(x_l)) \times n$

*important! for continuous distribution*

Since most distributions extend to infinity ( $\pm$ )

Careful when selecting own classes (bins) - avoid empty bins. Best practice **EZS (combine classes to form)**

eg1.

100 random values of  $x$  are taken from a distribution  $X$  which can take values 0, 1, 2, 3, 4 & 5, shown in the table. Test, at 1% significance level, whether  $X \sim B(5, 0.1)$  is a good fit.

|           |    |    |    |   |   |   |
|-----------|----|----|----|---|---|---|
| $x$       | 0  | 1  | 2  | 3 | 4 | 5 |
| Frequency | 60 | 21 | 15 | 2 | 1 | 1 |

*This is observed frequency (data sample)*

1. Null Hypothesis  $H_0: X \sim B(5, 0.1)$  is a good fit
2. Alternative Hypothesis  $H_1: X \sim B(5, 0.1)$  is NOT a good fit

$$3. \text{ Test statistic } C_n = \sum_{i=1}^l \frac{(O_i - E_i)^2}{E_i}$$

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

To find expected frequency,  $E_i$ :  $E_i = n \times P(X=x_i)$

$$P(X=x) = {}^x C_0 \times p^x \times (1-p)^{n-x}$$

$$E_0 = 100 P(X=0) = 100 ({}^5 C_0 \times 0.1^0 \times 0.9^5) = 59.049$$

$$E_1 = 100 P(X=1) = 100 ({}^5 C_1 \times 0.1^1 \times 0.9^4) = 32.805$$

$$\text{and so on... } E_2 = 7.29, E_3 = 0.81, E_4 = 0.045, E_5 = 0.001$$

| $x$   | 0      | 1      | 2    | 3    | 4     | 5     |
|-------|--------|--------|------|------|-------|-------|
| $O_i$ | 60     | 21     | 15   | 2    | 1     | 1     |
| $E_i$ | 59.049 | 32.805 | 7.29 | 0.81 | 0.045 | 0.001 |

these three  $E_i < S$ , we need to combine till no  $E_i < S$

so combine 2,3,4,5 (combining 3,4,5 won't make  $E_i \geq S$ !)

| $x$   | 0      | 1      | 2 to 5 |
|-------|--------|--------|--------|
| $O_i$ | 60     | 21     | 19     |
| $E_i$ | 59.049 | 32.805 | 8.146  |

$$C_n = \sum_{i=1}^l \frac{(O_i - E_i)^2}{E_i} = \frac{(60 - 59.049)^2}{59.049} + \frac{(21 - 32.805)^2}{32.805} + \frac{(19 - 8.146)^2}{8.146}$$

$$\therefore C_n = 18.7256$$

4. Distribution of test statistic  $C_n \sim \chi^2(l - 1 - k)$

5. Critical Region  $\mathcal{R}$   $C_n > \chi^2_{l-1-k, 1-\alpha}$  (e.g.  $\alpha = 1\%$ )

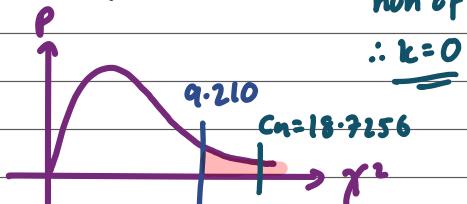
$$C_n \sim \chi^2(3 - 1 - 0)$$

$$C_n \sim \chi^2(2)$$

$k$  is number of estimated parameter.

In this question, we were told  $X \sim B(n, p)$   
none of the "n" and "p" are estimated.

$$\therefore \underline{k=0}$$



CHI-SQUARED Table - 2/2

|   |       |       |       |       |       |       |        |        |        |        |        |
|---|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|--------|
| 1 | 0.708 | 0.936 | 1.323 | 1.642 | 2.354 | 2.706 | 3.841  | 5.024  | 6.235  | 7.879  | 10.828 |
| 2 | 1.833 | 2.197 | 2.773 | 3.219 | 4.159 | 4.605 | 5.991  | 7.38   | 9.210  | 10.597 | 13.816 |
| 3 | 2.946 | 3.405 | 4.108 | 4.642 | 5.739 | 6.251 | 7.815  | 9.348  | 11.546 | 12.838 | 16.266 |
| 4 | 4.045 | 4.579 | 5.385 | 5.989 | 7.214 | 7.779 | 9.488  | 11.143 | 13.277 | 14.860 | 18.467 |
| 5 | 5.132 | 5.730 | 6.626 | 7.289 | 8.625 | 9.236 | 11.070 | 12.833 | 15.086 | 16.750 | 20.515 |

$$P(\chi^2 > \chi^2_{2, 0.01^2}) = 0.01$$

$$P(\chi^2 < \chi^2_{2, 0.01^2}) = 0.99$$

$$\chi^2_{2, 0.01^2} = 9.210$$

since  $C_n = 18.7256 > \chi^2_{2, 0.01^2} = 9.210$  (reject  $H_0$ ,  $X \sim B(5, 0.1)$  is NOT a good fit)

## eg2.

150 values of  $x$  are taken from a continuous population, shown in the table. Test whether  $N(\mu, \sigma^2)$  is a good fit at 10% significance level.

| $x$     | $f$ |
|---------|-----|
| -19     | 0   |
| 20 - 29 | 22  |
| 30 - 34 | 54  |
| 35 - 44 | 65  |
| 45 - 49 | 9   |
| 50 -    | 0   |

there is one unknown parameter

$$C_n \sim \chi^2(k - 1), k = 1$$

## Step 1. Estimate the unknown parameter

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

mean, so find mean of the sample.

$$\mu \approx \bar{x} = \frac{22(24.5) + 54(32) + 65(39.5) + 9(47)}{150} = 35.05$$

## Step 2. Define $H_0$ and $H_1$ .

$H_0: X \sim N(35.05, \sigma^2)$  is a good fit

$H_1: X \sim N(35.05, \sigma^2)$  is NOT a good fit

## Step 3. Find expected frequency, $E_i$

| $x$     | $f$ |
|---------|-----|
| -19     | 0   |
| 20 - 29 | 22  |
| 30 - 34 | 54  |
| 35 - 44 | 65  |
| 45 - 49 | 9   |
| 50 -    | 0   |

$$\begin{aligned} E_i &= 150 P(z_i \leq X \leq z_{i+1}) \\ &= 150 P(X \leq 19.5) \\ &= 150 P(19.5 \leq X \leq 29.5) \\ &= 150 P(29.5 \leq X \leq 34.5) \\ &= 150 P(34.5 \leq X \leq 44.5) \\ &= 150 P(44.5 \leq X \leq 49.5) \\ &= 150 P(X \geq 49.5) \end{aligned}$$

example how to find  $P(19.5 \leq X \leq 29.5)$

$$\begin{aligned} &= P(X \leq 29.5) - P(X \leq 19.5) \\ &= P\left(Z \leq \frac{29.5 - 35.05}{\sqrt{80}}\right) - P\left(Z \leq \frac{19.5 - 35.05}{\sqrt{80}}\right) \\ &= P(Z \leq -0.878) - P(Z \leq -2.459) \\ &= (1 - P(Z \leq 0.878)) - (1 - P(Z \leq 2.459)) \\ &= P(Z \leq 2.459) - P(Z \leq 0.878) \\ &= 0.993 - 0.8106 \\ &= 0.1825 \end{aligned}$$

## Step 4. Combine classes such that $E_i \geq 5$

| $x$     | $f$ | $E_i$  |
|---------|-----|--------|
| -19     | 0   | 1.046  |
| 20 - 29 | 22  | 27.469 |
| 30 - 34 | 54  | 41.288 |
| 35 - 44 | 65  | 70.063 |
| 45 - 49 | 9   | 8.466  |
| 50 -    | 0   | 1.675  |



| $x$     | $O_i$ | $E_i$  |
|---------|-------|--------|
| -19     | 0     | 22     |
| 20 - 29 | 22    | 28.515 |
| 30 - 34 | 54    | 41.288 |
| 35 - 44 | 65    | 70.063 |
| 45 - 49 | 9     | 10.141 |
| 50 -    | 0     |        |

$$\text{Step 5. Find } C_n = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

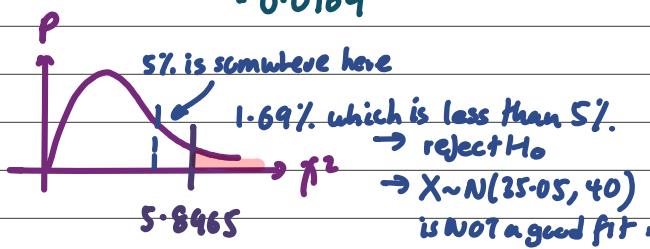
$$C_n = \frac{(22-28.515)^2}{28.515} + \frac{(54-41.288)^2}{41.288} + \frac{(65-70.063)^2}{70.063} + \frac{(9-10.141)^2}{10.141} = 5.8465$$

Step 6. Define test statistic and find critical region.

$$C_n \sim \chi^2(l-l-k) = \chi^2(4-1-1) = \chi^2(2)$$

(last example I find critical region, this question I try another approach — find p-value)

$$\begin{aligned} P(\chi^2 \geq 5.8965) &= 1 - P(\chi^2 < 5.8965) \\ &= 1 - 0.983 \\ &= 0.0169 \end{aligned}$$



CHI-SQUARED Table - 2/2

| $\nu$ | 60.0% | 66.7% | 75.0% | 80.0% | 87.5% | 90.0%  | 95.0%  | 97.5%  | 99.0%  | 99.5%  | 99.9%  |
|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|--------|--------|
| 1     | 0.708 | 0.936 | 1.323 | 1.642 | 2.354 | 2.706  | 3.841  | 5.024  | 6.635  | 8.79   | 10.828 |
| 2     | 1.833 | 2.197 | 2.773 | 3.219 | 4.159 | 4.605  | 5.991  | 7.378  | 9.210  | 10.597 | 13.816 |
| 3     | 2.946 | 3.405 | 4.108 | 4.642 | 5.739 | 6.251  | 7.815  | 9.348  | 11.345 | 12.838 | 16.266 |
| 4     | 4.045 | 4.579 | 5.385 | 5.989 | 7.214 | 7.779  | 9.488  | 11.143 | 13.277 | 14.860 | 18.467 |
| 5     | 5.132 | 5.730 | 6.626 | 7.289 | 8.625 | 9.236  | 11.070 | 12.833 | 15.086 | 16.750 | 20.515 |
| 6     | 6.211 | 6.867 | 7.841 | 8.558 | 0.009 | 10.645 | 12.509 | 14.440 | 16.819 | 18.518 | 22.458 |

$$0.975 + \frac{5.8965 - 5.024}{6.635 - 5.024} \times (0.99 - 0.975) \approx 0.983$$

ONE THING I NOTICE IS FOR YEAR 2'S STATS, IF  $E_i < 5$  WE DON'T COMBINE CLASSES. WE ALSO ALWAYS MINUS NUMBER OF PARAMETER FOR  $C_n \sim \chi^2(l-l-k)$  NO MATTER IF THEY ARE ESTIMATED OR NOT!

#### 8.4.2 Kolmogorov-Smirnov (KS) Test

- |                                   |                                                                  |
|-----------------------------------|------------------------------------------------------------------|
| 1. Null Hypothesis                | $H_0$ : Distribution is a good fit                               |
| 2. Alternative Hypothesis         | $H_1$ : Distribution is not a good fit                           |
| 3. Test statistic                 | $D_n = \max  F_n(x) - F_X(x) $                                   |
| 4. Distribution of test statistic | $\lim_{n \rightarrow \infty} P(\sqrt{n}D_n \leq z)$ is tabulated |
| 5. Critical Region $\mathcal{R}$  | $D_n > d_{n,1-\alpha}$ (e.g. $\alpha = 5\%$ )                    |
| 6. Evaluate $D_n$ under $H_0$     | if $d_{n0} > d_{n,1-\alpha} \Rightarrow$ Reject $H_0$            |
| 7. $p$ -value                     | $p(D_n > d_{n0})$                                                |

same eg from 8.4.1

150 values of  $x$  are taken from a continuous population, shown in the table. Test whether  $N(?, 40)$  is a good fit at 10 % significance level.

| $x$     | $f$ |
|---------|-----|
| -19     | 0   |
| 20 - 29 | 22  |
| 30 - 34 | 54  |
| 35 - 44 | 65  |
| 45 - 49 | 9   |
| 50 -    | 0   |

Step 1. Define  $H_0$  and  $H_1$ .

$H_0 : F(x) = F_T(x)$  for ALL values of  $x$

$H_1 : F(x) \neq F_T(x)$  for at least one value of  $x$

this is how we will  $H_0$  and  $H_1$  for ks-test!  
we are comparing empirical CDF to theoretical CDF

the reason for "one" value is because we only compare the max difference.

## Step 2. Find ECDF and Theoretical CDF at boundaries.

| Class Interval | Upper Boundary | Freq | Cum. Obs. Freq | ECDF   | Z-score | Theo. CDF | D      |
|----------------|----------------|------|----------------|--------|---------|-----------|--------|
| ≤ 19           | 19.5           | 0    | 0              | 0.0000 | -2.459  | 0.0070    | 0.0070 |
| 20 – 29        | 29.5           | 22   | 22             | 0.1467 | -0.878  | 0.1899    | 0.0432 |
| 30 – 34        | 34.5           | 54   | 76             | 0.5067 | -0.0875 | 0.4651    | 0.0416 |
| 35 – 44        | 44.5           | 65   | 141            | 0.9400 | 1.494   | 0.9324    | 0.0076 |
| 45 – 49        | 49.5           | 9    | 150            | 1.0000 | 2.285   | 0.9888    | 0.0112 |
| ≥ 50           | +∞             | 0    | 150            | 1.0000 | +∞      | 1.0000    | 0.0000 |

$$z = \frac{\text{upper boundary} - \mu}{\sigma}$$

$$= \frac{44.5 - 35.05}{\sqrt{40}}$$

Theoretical CDF =  $\Phi(z)$   
e.g. for  $z = 1.494$

NORMAL CUMULATIVE DISTRIBUTION FUNCTION

| x   | 0.00   | 0.01   | 0.02   | 0.03   | 0.04   | 0.05   | 0.06   | 0.07   | 0.08   | 0.09   |
|-----|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.0 | 0.5000 | 0.5040 | 0.5080 | 0.5120 | 0.5160 | 0.5199 | 0.5239 | 0.5279 | 0.5319 | 0.5359 |
| 0.1 | 0.5398 | 0.5438 | 0.5478 | 0.5517 | 0.5557 | 0.5596 | 0.5636 | 0.5675 | 0.5714 | 0.5753 |
| 0.2 | 0.5793 | 0.5832 | 0.5871 | 0.5910 | 0.5948 | 0.5987 | 0.6026 | 0.6064 | 0.6103 | 0.6141 |
| 0.3 | 0.6179 | 0.6217 | 0.6255 | 0.6293 | 0.6331 | 0.6368 | 0.6406 | 0.6443 | 0.6480 | 0.6517 |
| 0.4 | 0.6554 | 0.6591 | 0.6628 | 0.6664 | 0.6700 | 0.6736 | 0.6772 | 0.6808 | 0.6844 | 0.6879 |
| 0.5 | 0.6915 | 0.6950 | 0.6985 | 0.7019 | 0.7054 | 0.7088 | 0.7123 | 0.7157 | 0.7190 | 0.7224 |
| 0.6 | 0.7270 | 0.7298 | 0.7327 | 0.7357 | 0.7385 | 0.7414 | 0.7444 | 0.7474 | 0.7503 | 0.7531 |
| 0.7 | 0.7626 | 0.7651 | 0.7674 | 0.7702 | 0.7730 | 0.7754 | 0.7781 | 0.7809 | 0.7832 | 0.7852 |
| 0.8 | 0.7981 | 0.7910 | 0.7939 | 0.7967 | 0.7995 | 0.8023 | 0.8051 | 0.8078 | 0.8109 | 0.8133 |
| 0.9 | 0.8150 | 0.8186 | 0.8212 | 0.8238 | 0.8264 | 0.8290 | 0.8315 | 0.8340 | 0.8365 | 0.8389 |
| 1.0 | 0.8413 | 0.8438 | 0.8461 | 0.8485 | 0.8508 | 0.8531 | 0.8554 | 0.8577 | 0.8599 | 0.8621 |
| 1.1 | 0.8643 | 0.8665 | 0.8686 | 0.8708 | 0.8729 | 0.8750 | 0.8770 | 0.8790 | 0.8810 | 0.8830 |
| 1.2 | 0.8849 | 0.8869 | 0.8888 | 0.8907 | 0.8925 | 0.8944 | 0.8962 | 0.8980 | 0.8997 | 0.9015 |
| 1.3 | 0.9032 | 0.9049 | 0.9066 | 0.9082 | 0.9099 | 0.9115 | 0.9131 | 0.9147 | 0.9162 | 0.9177 |
| 1.4 | 0.9192 | 0.9207 | 0.9222 | 0.9236 | 0.9251 | 0.9265 | 0.9279 | 0.9292 | 0.9305 | 0.9319 |
| 1.5 | 0.9332 | 0.9345 | 0.9357 | 0.9370 | 0.9382 | 0.9394 | 0.9406 | 0.9418 | 0.9429 | 0.9441 |

## Step 3. Define test statistic

$$D_n = \max(D) = 0.0432$$

Kolmogorov-Smirnov Test (level of significance  $\alpha$ )

| n    | 0.2       | 0.1       | 0.05      | 0.02      | 0.01      |
|------|-----------|-----------|-----------|-----------|-----------|
| 1    | 0.9000    | 0.9500    | 0.9750    | 0.9900    | 0.9950    |
| 2    | 0.6838    | 0.7764    | 0.8419    | 0.9000    | 0.9293    |
| 3    | 0.5648    | 0.6360    | 0.7076    | 0.7846    | 0.8290    |
| 4    | 0.4927    | 0.5652    | 0.6239    | 0.6889    | 0.7342    |
| 5    | 0.4470    | 0.5094    | 0.5633    | 0.6272    | 0.6685    |
| 6    | 0.4104    | 0.4680    | 0.5193    | 0.5774    | 0.6166    |
| 7    | 0.3815    | 0.4361    | 0.4834    | 0.5384    | 0.5758    |
| 8    | 0.3583    | 0.4096    | 0.4543    | 0.5065    | 0.5418    |
| 9    | 0.3391    | 0.3875    | 0.4300    | 0.4796    | 0.5133    |
| 10   | 0.3226    | 0.3687    | 0.4092    | 0.4566    | 0.4889    |
| 11   | 0.3083    | 0.3524    | 0.3912    | 0.4367    | 0.4677    |
| 12   | 0.2958    | 0.3382    | 0.3754    | 0.4192    | 0.4490    |
| 13   | 0.2847    | 0.3255    | 0.3614    | 0.4036    | 0.4325    |
| 14   | 0.2748    | 0.3142    | 0.3489    | 0.3897    | 0.4176    |
| 15   | 0.2659    | 0.3040    | 0.3376    | 0.3771    | 0.4042    |
| 16   | 0.2578    | 0.2947    | 0.3273    | 0.3657    | 0.3920    |
| 17   | 0.2504    | 0.2863    | 0.3180    | 0.3553    | 0.3809    |
| 18   | 0.2436    | 0.2785    | 0.3094    | 0.3457    | 0.3706    |
| 19   | 0.2373    | 0.2714    | 0.3014    | 0.3369    | 0.3612    |
| 20   | 0.2316    | 0.2647    | 0.2941    | 0.3287    | 0.3524    |
| 21   | 0.2262    | 0.2586    | 0.2872    | 0.3210    | 0.3443    |
| 22   | 0.2212    | 0.2528    | 0.2809    | 0.3139    | 0.3367    |
| 23   | 0.2165    | 0.2475    | 0.2749    | 0.3073    | 0.3295    |
| 24   | 0.2120    | 0.2424    | 0.2693    | 0.3010    | 0.3229    |
| 25   | 0.2079    | 0.2377    | 0.2640    | 0.2952    | 0.3166    |
| 26   | 0.2040    | 0.2332    | 0.2591    | 0.2896    | 0.3106    |
| 27   | 0.2003    | 0.2290    | 0.2544    | 0.2844    | 0.3050    |
| 28   | 0.1968    | 0.2250    | 0.2499    | 0.2794    | 0.2997    |
| 29   | 0.1935    | 0.2212    | 0.2457    | 0.2747    | 0.2947    |
| 30   | 0.1903    | 0.2176    | 0.2417    | 0.2702    | 0.2899    |
| 31   | 0.1873    | 0.2141    | 0.2379    | 0.2660    | 0.2853    |
| 32   | 0.1844    | 0.2108    | 0.2342    | 0.2619    | 0.2809    |
| 33   | 0.1817    | 0.2077    | 0.2308    | 0.2580    | 0.2768    |
| 34   | 0.1791    | 0.2047    | 0.2274    | 0.2543    | 0.2728    |
| 35   | 0.1766    | 0.2018    | 0.2242    | 0.2507    | 0.2690    |
| 36   | 0.1742    | 0.1991    | 0.2212    | 0.2473    | 0.2653    |
| 37   | 0.1719    | 0.1965    | 0.2183    | 0.2440    | 0.2618    |
| 38   | 0.1697    | 0.1939    | 0.2154    | 0.2409    | 0.2584    |
| 39   | 0.1675    | 0.1915    | 0.2127    | 0.2379    | 0.2552    |
| 40   | 0.1655    | 0.1893    | 0.2101    | 0.2349    | 0.2521    |
| > 40 | 1.07 / √n | 1.22 / √n | 1.36 / √n | 1.52 / √n | 1.63 / √n |

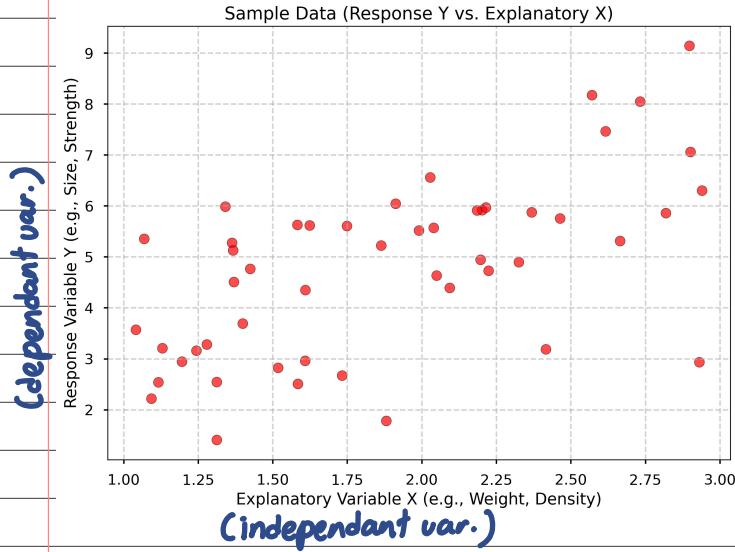
$$d_{150, 0.9} = \frac{1.22}{\sqrt{150}} = 0.0996$$

## Step 4. Conclude

$D_n = 0.0432 < d_{150, 0.9} = 0.0996$  (do not reject  $H_0 \rightarrow X \sim N(35.05, 40)$  is a good fit)

## C9. Linear Regression

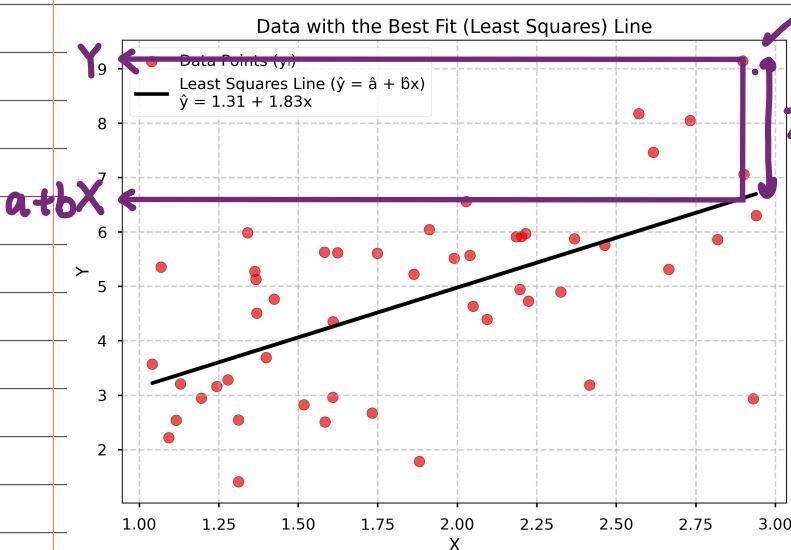
### 9.1 Introduction



Linear regression is to fit a straight line onto a dataset of dependant vs. independant  
The equation of Linear Regression is :

$$Y = a + bX + Z \quad \begin{array}{l} \text{Z is an error term: 1. Z is an independent term.} \\ \text{(residual)} \end{array}$$

a and b are called regression coefficient.



this is only "Z" of one of the point from the samples.  
For n points in a samples, there is a number of X, Y and Z  
cool thing to note is for all n number of Z, the mean of all the Z is 0!  
 $\rightarrow E[Z] = 0$ !  
this infer that the mean of all the Y is  $a+bX$ !  
 $\rightarrow E[Y] = a+bX$

about a and b:

in real world scenario, we can't measure population, but we can obtain sample to estimate population ...

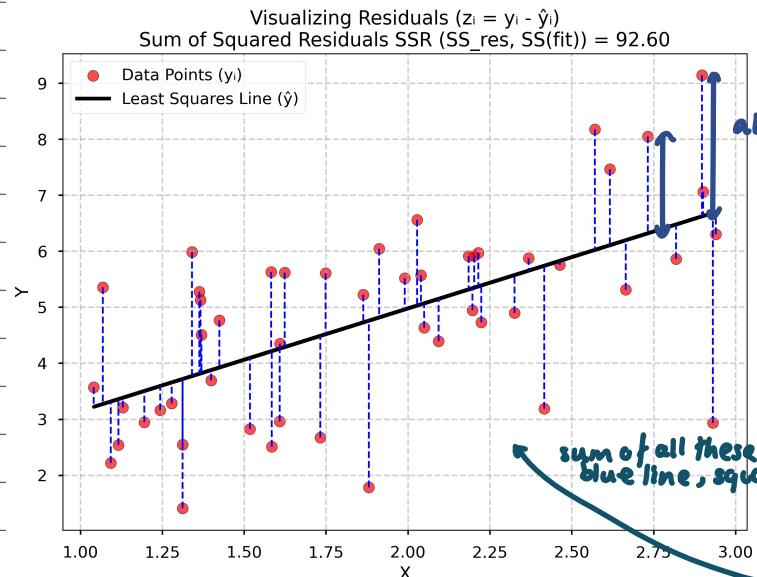
eg. we can't find the mean of human height (there is 8 billion human  $\rightarrow \mu$ )

but we can obtain a sample and get the sample mean  $\rightarrow \bar{z}$

same for a and b. We can't get the entire population data to fit a linear line, but we can obtain a sample and fit a linear line.

Hence, the "a" and "b" we get for our line is named " $\hat{a}$ " and " $\hat{b}$ "  
(same idea as " $\bar{x}$ " for " $\mu$ " and " $S_x^2$ " for " $O_x^2$ ")

So how do we obtain  $\hat{a}$  and  $\hat{b}$  that will give the best fit?



all these blue lines are  $z$  (residual)  
the idea is: line of best fit should have small residual.

Hence: we will try to obtain

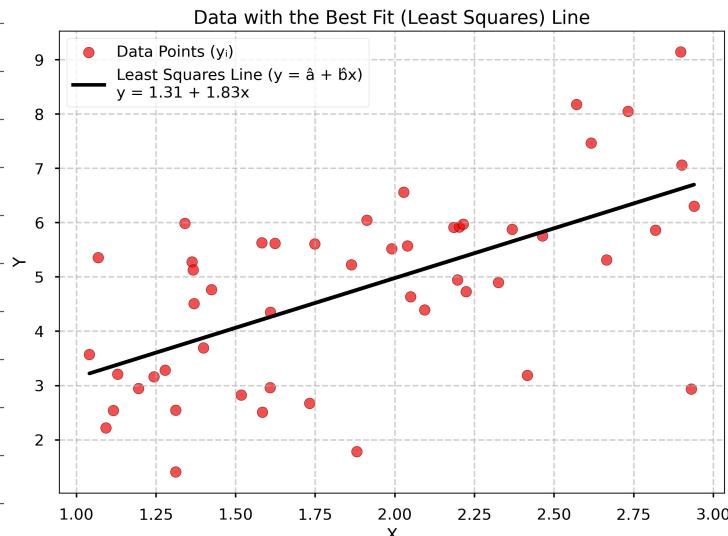
$$\min \left\{ \sum_{i=1}^N [y_i - (\hat{a} + \hat{b}x_i)]^2 \right\}$$

$y_i - (\hat{a} + \hat{b}x_i) = z_i$   
cause  $y = \hat{a} + \hat{b}x + z$

we call this sum of squared residual  
cause  $\min \left\{ \sum z_i^2 \right\}$

For any sets of sample, like the one shown in the graph above, will have ONE sets of  $\hat{a}$  and  $\hat{b}$  that gives line of best fits.

In this example,  $\hat{a} = 1.31$  and  $\hat{b} = 1.83$ :



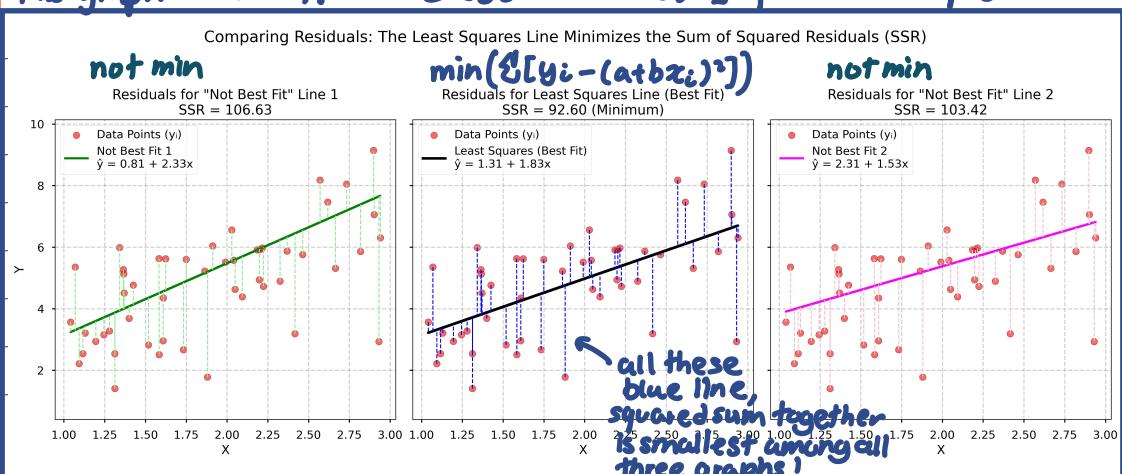
\* There is a shortcut in solving  
 $\min \left\{ \sum_{i=1}^N [y_i - (\hat{a} + \hat{b}x_i)]^2 \right\}$

instead of plotting hundreds of graph and find the best  $\hat{a}$  and  $\hat{b}$ :

$$\hat{b} = \frac{\text{cov}_{x,y}}{\text{var}_x}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

↑  
sample covariance  
↑  
sample mean y  
↓  
sample variance x  
↑  
sample mean z

this graph shows different  $\sum [y_i - (\hat{a} + \hat{b}x_i)]^2$  for same sample:



$\min \left\{ \sum_{i=1}^n [y_i - (\hat{a} + \hat{b}x_i)]^2 \right\}$  is only for determining  $\hat{a}$  and  $\hat{b}$  for the line of best fit. However, it doesn't tell us how good of a fit, compare to other samples.

The goodness-of fit of the regression model is measured with the **coefficient of determination**  $r^2$ :

$$r^2 = \frac{\text{var}_{\hat{a}+\hat{b}x}}{\text{var}_y} = c_{xy}^2 \leq 1$$

$$r^2 = C_{xy}^2$$

or

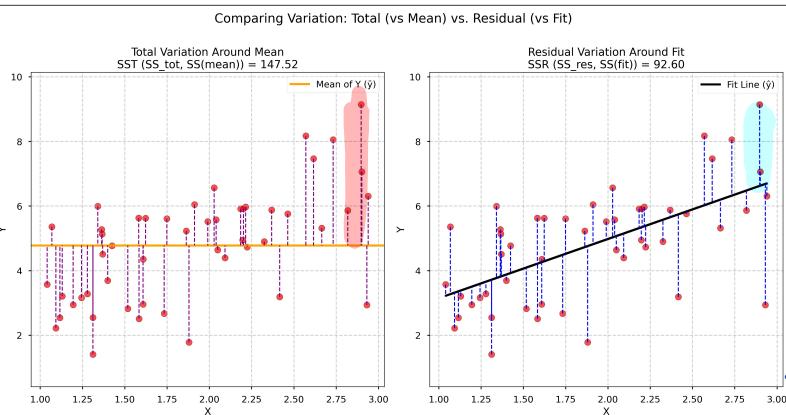
$$r^2 = \frac{\text{cov}_{xy}^2}{\text{Var}_x \text{Var}_y}$$

Don't confuse with correlation coefficient ( $C_{xy} = \frac{\text{cov}_{xy}}{S_x S_y}$ )

$$r^2 = \frac{\sum_{i=1}^n (\hat{a} + \hat{b}x_i - (\bar{a} + \bar{b}\bar{x}))^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \left( \frac{\text{cov}_{xy}}{\text{var}_x} \right)^2 \frac{\text{var}_x}{\text{var}_y} = \frac{\text{cov}_{xy}^2}{\text{var}_x \text{var}_y}$$

JUST  
REMEMBER  
THE FORMULA

THIS IS NOT  
IMPORTANT  
DON'T NEED  
UNDERSTAND



this is how to visualise  $r^2$ :  
it's the blue line distance minus  
its corresponding red line distance .

Basically, population parameter (e.g.  $\mu$ ,  $\sigma^2$ ) are unknown. To estimate we first find (e.g.  $\bar{x}$  and  $s^2$ ). Then referring to c7, we can find CI of  $\mu$  and  $\sigma^2$ .

## 9.2 Properties of Parameter Estimates

Remember in c7, we know we can find  $\bar{x}$ , then estimate for  $\mu$ :

| 7.4.1 Sampling Distribution of Sample Mean.                                                                                          |                                                                                                                                                                                 |
|--------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Sample Mean                                                                                                                          | basically we will try to use $n\bar{X}$ . If $\sigma^2$ is known, if not known, try to estimate it. Only if not possible to estimate ( $n$ small) we will use t-test.           |
| Standard Deviation (Population) is known                                                                                             | Standard Deviation (Population) is not known.                                                                                                                                   |
| → if was told $X$ is normally distributed, no matter $n$ is big ( $>30$ ) or small, $\bar{X} \sim N(\mu, \sigma^2/n)$                | → if $n$ is small ( $n < 30$ ) or was told not to estimate $\sigma^2$ $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim t(n-1)$ <small>(t-distribution? check example 3 below!)</small> |
| → if $n$ is big ( $>30$ ), $\bar{X} \sim N(\mu, \sigma^2/n)$                                                                         | → if $n$ is big ( $n > 30$ ), and is possible to estimate $\sigma^2$ (given dist. of $X$ ), $\bar{X} \sim N(\mu, \sigma^2/n)$ with estimated $\sigma^2$                         |
| for year 2 - statistics, don't worry about:<br>- s.d. known, $X$ is not normally dist. AND $n$ is small. (technically we use t-test) |                                                                                                                                                                                 |

the estimates is just an interval, e.g. 95% confidence that  $\mu$  is between (5, 10) given  $\bar{x}$  is 7.5

but to do this estimation, we need to know:

1. What is the distribution of the sample statistic: e.g.  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$
2. we need to find the parameter of the distribution: e.g.  $\mu$  and  $\frac{\sigma^2}{n}$

In c9, we are estimating for  $a$  and  $b$  with  $\hat{a}$  and  $\hat{b}$

From 9.1 Introduction, we know we can find  $\hat{a}$  and  $\hat{b}$  but not  $a$  and  $b$ :

$$\hat{b} = \frac{\text{cov}_{x,y}}{\text{var}_x} \quad \text{and} \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

Not going into too much details (cause it's very complicated), you just need to know that:

$$\hat{a} \sim N(E(\hat{a}), \text{var}(\hat{a})) ; \quad \hat{b} \sim N(E(\hat{b}), \text{var}(\hat{b}))$$

So what we need to find?  $E(\hat{a})$ ,  $\text{var}(\hat{a})$ ,  $E(\hat{b})$ ,  $\text{var}(\hat{b})$



$$E(\hat{a}) = a \quad \text{and} \quad E(\hat{b}) = b$$

} we can say that  $\hat{a}$  and  $\hat{b}$  are unbiased estimator of  $a$  and  $b$ . same like  $\bar{x}$  because  $E(\bar{x}) = \mu$ .

$$\text{Var}(\hat{a}) = \left( \frac{1}{n} + \frac{\bar{x}^2}{n\text{var}_x} \right) \text{Var}(Z), \quad \text{and} \quad \text{Var}(\hat{b}) = \frac{1}{n\text{var}_x} \text{Var}(Z)$$

However,  $\text{Var}(Z)$  itself is a POPULATION PARAMETER that needs to be estimated:

$$\text{Var}(Z) \approx \widehat{\text{Var}(Z)} = \frac{1}{n-2} \left( n\text{var}_y - n \frac{\text{cov}_{xy}^2}{\text{var}_x} \right)$$

this is mathematically equivalent to:

$$\sum [y_i - (\hat{a} + \hat{b}z_i)]^2$$

because of this  $n-2$ , the  $\text{var}_x$ ,  $\text{var}_y$ , and  $\text{cov}_{xy}$  should be population (biased)

small info on what is  $\text{Var}(Z)$  ↗ population parameter  
 $\text{Var}(Z)$  is 'true error variance' which is the sum of square residual divided by  $n$ , where  $n$  is size of population.  
 But since it is population, we can only estimate it with  $\text{Var}(\hat{Z})$  = sum of square residual divided by  $n-2$ , where  $n$  is size of sample.  
 So  $\text{Var}(Z)$  is actually very related to the residual (the blue line on the graph shown just now)

### 9.3 Confidence Interval (Population Parameters: a and b Estimate)

$$\hat{a} \sim N(E(\hat{a}), \text{Var}(\hat{a})), \quad \hat{b} \sim N(E(\hat{b}), \text{Var}(\hat{b}))$$

$$\hat{a} \sim N\left(a, \left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right) \text{Var}(Z)\right), \quad \hat{b} \sim N\left(b, \frac{\text{Var}(Z)}{n\text{Var}_x}\right).$$

both of these requires population statistic of  $\text{Var}(Z)$ ,

Two possible scenario :

1. If question straight away give us  $\text{Var}(Z)$ : USE NORMAL DISTRIBUTION

$$\frac{\hat{X} - E(X)}{\sqrt{\text{Var}(X)}} \sim N(0,1)$$

$\star$  straight away use the given  $\text{Var}(Z)$ !

$$P\left(-z_{\alpha/2} < \frac{\hat{a} - a}{\sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right) \text{Var}(Z)}} < z_{\alpha/2}\right) = 1 - \alpha$$

$$P\left(\hat{a} - z_{\alpha/2} \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right) \text{Var}(Z)} < a < \hat{a} + z_{\alpha/2} \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right) \text{Var}(Z)}\right) = 1 - \alpha$$

$$\frac{\hat{b} - b}{\sqrt{\text{Var}(Z)/n\text{Var}_x}} \sim N(0,1)$$

$$P\left(-z_{\alpha/2} < \frac{\hat{b} - b}{\sqrt{\text{Var}(Z)/n\text{Var}_x}} < z_{\alpha/2}\right) = 1 - \alpha$$

$$P\left(\hat{b} - z_{\alpha/2} \sqrt{\frac{\text{Var}(Z)}{n\text{Var}_x}} < b < \hat{b} + z_{\alpha/2} \sqrt{\frac{\text{Var}(Z)}{n\text{Var}_x}}\right) = 1 - \alpha$$

2. If question does not give us  $\text{Var}(Z)$ : USE T DISTRIBUTION (with estimated  $\text{Var}(Z)$ )

$$\widehat{\text{Var}(Z)} = \frac{1}{n-2} \left( n\text{Var}_y - n \frac{\text{cov}_{xy}^2}{\text{Var}_x} \right)$$

$\star$  remember: ALWAYS "n-2"!

$$\frac{\hat{a} - a}{\sqrt{\widehat{\text{Var}(Z)} \left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right)}} \sim t(n-2)$$

$\star$  Estimate  $\text{Var}(Z)$  with  $\text{Var}(\hat{Z})$  first!

$$P\left(-t_{n-2,\alpha/2} < \frac{\hat{a} - a}{\sqrt{\widehat{\text{Var}(Z)} \left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right)}} < t_{n-2,\alpha/2}\right) = 1 - \alpha$$

$$P\left(\hat{a} - t_{n-2,\alpha/2} \sqrt{\widehat{\text{Var}(Z)} \left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right)} < a < \hat{a} + t_{n-2,\alpha/2} \sqrt{\widehat{\text{Var}(Z)} \left(\frac{1}{n} + \frac{\bar{x}^2}{n\text{Var}_x}\right)}\right) = 1 - \alpha$$

$$\frac{\hat{b} - b}{\sqrt{\widehat{\text{Var}(Z)}/n\text{Var}_x}} \sim t(n-2)$$

$$P\left(-t_{n-2,\alpha/2} < \frac{\hat{b} - b}{\sqrt{\widehat{\text{Var}(Z)}/n\text{Var}_x}} < t_{n-2,\alpha/2}\right) = 1 - \alpha$$

$$P\left(\hat{b} - t_{n-2,\alpha/2} \sqrt{\frac{\widehat{\text{Var}(Z)}}{n\text{Var}_x}} < b < \hat{b} + t_{n-2,\alpha/2} \sqrt{\frac{\widehat{\text{Var}(Z)}}{n\text{Var}_x}}\right) = 1 - \alpha$$

(An example for this 9.3 will be given at the very end of this chapter)

When told to find an estimate value for  $Y$  for a given  $X$ , should find both  $E(Y|X=x_0)$  and  $Y(X=x_0)$ 's interval!

## 9.4 Confidence Interval for $E(Y|X=x_0)$ and $Y_{X=x_0}$

### 1. Confidence Interval for the Mean Value of $Y$ at $X = x_0$

- Question: If we observed many, many  $Y$  values when  $X$  is exactly  $x_0$ , what is the likely range for the average of all those  $Y$  values?

understand what are we finding!

$$\text{average of all those } Y \text{ values: } E(Y | X = x_0) = a + bx_0$$

$$\frac{(\hat{a} + \hat{b}x_0) - (a + bx_0)}{\sqrt{\text{Var}(Z) \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)}} \sim t(n-2)$$

$$P\left(-t_{n-2,\alpha/2} < \frac{(\hat{a} + \hat{b}x_0) - (a + bx_0)}{\sqrt{\text{Var}(Z) \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)}} < t_{n-2,\alpha/2}\right) = 1 - \alpha$$

$$P\left((\hat{a} + \hat{b}x_0) - t_{n-2,\alpha/2} \cdot \sqrt{\text{Var}(Z) \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)} < a + bx_0 < (\hat{a} + \hat{b}x_0) + t_{n-2,\alpha/2} \cdot \sqrt{\text{Var}(Z) \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)}\right) = 1 - \alpha$$

### 2. Prediction Interval for an Individual Value of $Y$ at $X = x_0$

same as confidence interval

- Question: If we observe one single new  $Y$  value when  $X$  is  $x_0$ , what is the likely range for that specific  $Y$  value?

$$\text{new } Y \text{ value when } X \text{ is } x_0: Y_{(X=x_0)} = a + bx_0 + Z$$

$$\frac{(a + bx_0 + Z) - (\hat{a} + \hat{b}x_0)}{\sqrt{\text{Var}(Z) \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)}} \sim t(n-2)$$

$$P\left(-t_{n-2,\alpha/2} < \frac{(a + bx_0 + Z) - (\hat{a} + \hat{b}x_0)}{\sqrt{\text{Var}(Z) \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)}} < t_{n-2,\alpha/2}\right) = 1 - \alpha$$

$$P\left((\hat{a} + \hat{b}x_0) - t_{n-2,\alpha/2} \cdot \sqrt{\text{Var}(Z) \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)} < (a + bx_0 + Z) < (\hat{a} + \hat{b}x_0) + t_{n-2,\alpha/2} \cdot \sqrt{\text{Var}(Z) \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{n \text{var}_x} \right)}\right) = 1 - \alpha$$

basically same steps as c8. Looks more complicated, but is just the same.

## 9.5 Statistical Testing (Hypothesis Testing)

Three different test will be examined

### (a) Test on the slope parameter $b$

Most commonly this is used to test whether there is a linear relation or not. Therefore, the *null* hypothesis will be  $H_0 : b = 0$ . Generally, we consider

$$H_0 : b = b^* \quad \text{and} \quad H_1 : b \neq b^*$$

( so commonly )  
 $H_0 : b=0$   
 $H_1 : b\neq 0$

### (b) Test on the correlation

The sample correlation coefficient can be used to decide whether two variables  $X$  and  $Y$  are linearly related (i.e. population  $\text{Corr}(X, Y) \neq 0$  ).

### (c) Test on the regression - *F*-test

This test assess the quality of the regression using the *F* distribution.

### 9.5.1 Slope Parameter Test

#### (a) Test on the slope parameter $b$

1. Null Hypothesis  $H_0 : b = b^*$
2. Alternative Hypothesis  $H_1 : b \neq b^*$
3. Test statistic  $T = \frac{\hat{b} - b}{\sqrt{\frac{\text{Var}(\hat{b})}{n \text{Var}_x}}}$
4. Distribution of test statistic  $T \sim t(n-2)$
5. Critical Region  $\mathfrak{R}$   $|T| > t_{n-2,\alpha/2}$
6. Evaluate  $T$  under  $H_0$  if  $|t_0| > t_{n-2,\alpha/2} \Rightarrow \text{Reject}$
7.  $p$ -value  $p(|T| > t_0)$

This can easily be transformed to one-sided test where the alternative hypothesis could be  $H_1 : b > b^*$

(Example for this question will be provided at the very end of this chapter)

## 9.5.2 Correlation Test (between X and Y)

### (b) Test on the correlation

1. Null Hypothesis  $H_0 : \text{Corr}(X, Y) = 0$
2. Alternative Hypothesis  $H_1 : \text{Corr}(X, Y) \neq 0$
3. Test statistic  $V = \frac{\sqrt{n-3}}{2} \ln \left[ \frac{(1+c_{xy})(1-\text{Corr}(X,Y))}{(1-c_{xy})(1+\text{Corr}(X,Y))} \right]$
4. Distribution of test statistic  $V \sim N(0, 1)$
5. Critical Region  $\mathfrak{R}$   $|V| > z_{\alpha/2}$
6. Evaluate  $V$  under  $H_0$  if  $|v_0| > z_{\alpha/2} \Rightarrow \text{Reject } H_0$
7.  $p - \text{value}$   $p(|V| > v_0)$

where  $c_{xy}$  is the sample correlation coefficient.

In practice, the correlation test does not require the use of the regression coefficients.

## 9.5.3 Regression Test (F-test)

### (c) Test on the regression - F-test

1. Null Hypothesis  $H_0 : b = 0$
2. Alternative Hypothesis  $H_1 : b \neq 0$
3. Test statistic  $F = \frac{ss_{reg}}{ss_{res}/n-2}$
4. Distribution of test statistic  $F \sim F(1, n - 2)$
5. Critical Region  $\mathfrak{R}$   $F > f_{1,n-2,\alpha}$
6. Evaluate  $F$  under  $H_0$  if  $f_0 > f_{1,n-2,\alpha} \Rightarrow \text{Reject } H_0$
7.  $p - \text{value}$   $p(F > f_0)$

The test statistic can also be written as:

$$F = \frac{\sum_{i=1}^n (\hat{a} + \hat{b}x_i - \bar{y})^2}{\sum_{i=1}^n (y_i - (\hat{a} + \hat{b}x_i))^2 / (n - 2)}$$

$$ss_{reg} = \sum_{i=1}^n (\hat{a} + \hat{b}x_i - (\hat{a} + \hat{b}\bar{x}))^2 = n \text{var}_{\hat{a} + \hat{b}x}$$

$$ss_{res} = \sum_{i=1}^n (y_i - (\hat{a} + \hat{b}x_i))^2 = (n - 2) \widehat{\text{Var}(\bar{Z})}$$

$$ss_{tot} = ss_{reg} + ss_{res} = \sum_{i=1}^n (y_i - \bar{y})^2 = n \text{var}_y$$

not important, but if you want to know what is f-distribution...

If  $U \sim \chi^2(v_1)$  and  $V \sim \chi^2(v_2)$

$$\frac{U/v_1}{V/v_2} \sim F(v_1, v_2)$$

f-distribution is ratio of two chi-squared distribution.

eg.

#### 4. Relation between applied force and length of a wire

Ten steel wires of diameter 0.5 mm and length 2.5 m were extended in a laboratory by applying vertical forces of varying magnitudes. The results are as follows:

|                      |     |     |     |     |     |     |     |     |     |     |
|----------------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Force (kg)           | 15  | 19  | 25  | 35  | 42  | 48  | 53  | 56  | 62  | 65  |
| Length increase (mm) | 1.7 | 2.1 | 2.5 | 3.4 | 3.9 | 4.9 | 5.4 | 5.7 | 6.6 | 7.2 |

a - We wish to explore the correlation between Force and Length increase. Test the hypothesis that the population correlation is equal to 0. (*Significant level  $\alpha$  is not given, let us use  $\alpha = 5\%$ .*)

#### (b) Test on the correlation

1. Null Hypothesis  $H_0 : \text{Corr}(X, Y) = 0$
2. Alternative Hypothesis  $H_1 : \text{Corr}(X, Y) \neq 0$
3. Test statistic  $V = \frac{\sqrt{n-3}}{2} \ln \left[ \frac{(1+c_{xy})(1-\text{Corr}(X,Y))}{(1-c_{xy})(1+\text{Corr}(X,Y))} \right]$
4. Distribution of test statistic  $V \sim N(0, 1)$
5. Critical Region  $\mathcal{R}$   $|V| > z_{\alpha/2}$
6. Evaluate  $V$  under  $H_0$  if  $|v_0| > z_{\alpha/2} \Rightarrow \text{Reject } H_0$
7.  $p$ -value  $p(|V| > v_0)$

**Step 1. Define  $H_0$  and  $H_1$ ,**

$$H_0 : \text{Corr}(X, Y) = 0, \quad H_1 : \text{Corr}(X, Y) \neq 0$$

**Step 2. Define Test Statistic**

3. Test statistic

$$V = \frac{\sqrt{n-3}}{2} \ln \left[ \frac{(1+c_{xy})(1-\text{Corr}(X,Y))}{(1-c_{xy})(1+\text{Corr}(X,Y))} \right]$$

(need to find sample's correlation coefficient)

| 1.2.4 Measure Relation between Datasets. (check whether there is any linear correlation between $x_i$ and $y_i$ ) |                                                                                                               |
|-------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------|
| 1. SAMPLE covariance.                                                                                             | (problem about this is covariance has units, if $x$ and $y$ have unit [m] then covariance has unit [ $m^2$ ]) |
| $\text{cov}_{x,y} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$                                            | 2. SAMPLE correlation coefficient.<br>$c_{xy} = \frac{\text{cov}_{x,y}}{S_x S_y} \quad -1 \leq c_{xy} \leq 1$ |
| → to non-dimensionalise...                                                                                        |                                                                                                               |

From the sample,

$$n=10$$

$$\bar{x} (\text{mean of Force}) = 42.0$$

$$\bar{y} (\text{mean of length increase}) = 4.34$$

$$\sum (x_i - \bar{x})(y_i - \bar{y}) = 305.7$$

$$\text{cov}_{x,y} = \frac{305.7}{10-1} = 33.9667$$

$$S_x = 17.8823 \quad (s \rightarrow \text{sample's s.s.d.})$$

$$S_y = 1.9156$$

$$c_{xy} = \frac{33.9667}{17.8823 \times 1.9156} = 0.9916$$

$$\text{Test statistic, } V = \frac{\sqrt{n-3}}{2} \ln \left( \frac{(1+0.9916)(1-0)}{(1-0.9916)(1+0)} \right)$$

$$V = \frac{\sqrt{n-3}}{2} \ln \left[ \frac{(1+c_{xy})(1-\text{Corr}(X,Y))}{(1-c_{xy})(1+\text{Corr}(X,Y))} \right]$$

#

$\text{Corr}(X,Y) = 0$   
cause at this stage, we assume  $H_0$  is true!

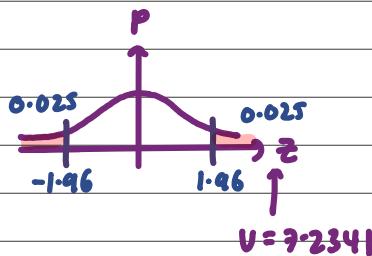
$$= 7.2341$$

### Step 3. Define critical region

$$P(Z > z_{\alpha/2}) = 0.025$$

$$P(Z < z_{\alpha/2}) = 0.975$$

$$z_{\alpha/2} = 1.96$$



$$\therefore U = 7.2341 > z_{\alpha/2} = 1.96$$

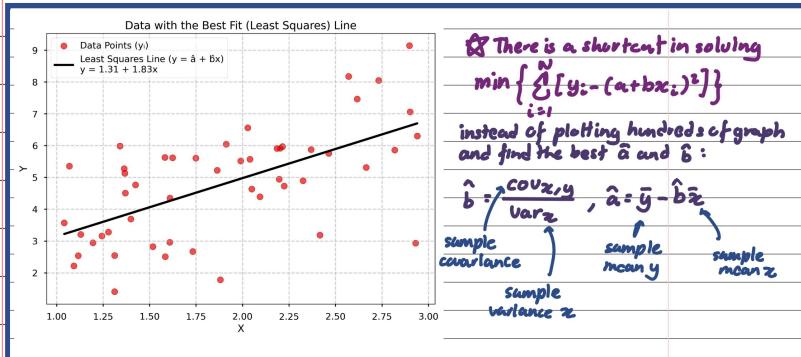
Reject  $H_0 \rightarrow \text{Corr}(X, Y) \neq 0$

(find  $\hat{a}$  and  $\hat{b}$ )

b - Estimate the parameters of a simple linear regression model with Force as an

explanatory variable, indicating the coefficient of determination  
independent variable (z) (find  $r^2$ )

To find  $\hat{a}$  and  $\hat{b}$ :



$$\hat{b} = \frac{\text{cov}_{xy}}{\text{var}_x} = \frac{33.9667}{319.9767} = 0.1062$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x} = 4.34 - 0.1062(42.0) = -0.1204$$

$$\left. \right\} \hat{y} = -0.1204 + 0.1062 \hat{x}$$

To find  $r^2$ :

$$r^2 = \frac{\sum_{i=1}^n (\hat{a} + \hat{b}x_i - (\bar{a} + \bar{b}\bar{x}))^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \left( \frac{\text{cov}_{xy}}{\text{var}_x} \right)^2 \frac{\text{var}_x}{\text{var}_y} = \frac{\text{cov}_{xy}^2}{\text{var}_x \text{var}_y}$$

$$r^2 = \frac{33.9667^2}{319.9767 \times 3.6695} = 0.9832$$

(This  $r^2$  result indicates 98.32% of Y can be explained with a linear relationship with X)

c - Find 95% confidence limits for the two parameters

## 2. If question does not give us $\text{Var}(Z)$ : USE T DISTRIBUTION (with estimated $\text{Var}(Z)$ )

$$\widehat{\text{Var}}(Z) = \frac{1}{n-2} \left( \text{near}_y - \frac{\text{cov}_{xy}^2}{\text{var}_x} \right)$$

*\*remember: ALWAYS "n-2"!*

**\* Estimate  $\text{Var}(Z)$  with  $\text{Var}(\hat{Z})$  first!**

$$\sqrt{\frac{\hat{a} - a}{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)}} \sim t(n-2)$$

$$P \left( -t_{n-2,\alpha/2} < \frac{\hat{a} - a}{\sqrt{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)}} < t_{n-2,\alpha/2} \right) = 1 - \alpha$$

$$P \left( \hat{a} - t_{n-2,\alpha/2} \sqrt{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)} < a < \hat{a} + t_{n-2,\alpha/2} \sqrt{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)} \right) = 1 - \alpha$$

$$\sqrt{\frac{\hat{b} - b}{\widehat{\text{Var}}(Z) / n \text{var}_x}} \sim t(n-2)$$

$$P \left( -t_{n-2,\alpha/2} < \frac{\hat{b} - b}{\sqrt{\widehat{\text{Var}}(Z) / n \text{var}_x}} < t_{n-2,\alpha/2} \right) = 1 - \alpha$$

$$P \left( \hat{b} - t_{n-2,\alpha/2} \sqrt{\frac{\widehat{\text{Var}}(Z)}{n \text{var}_x}} < b < \hat{b} + t_{n-2,\alpha/2} \sqrt{\frac{\widehat{\text{Var}}(Z)}{n \text{var}_x}} \right) = 1 - \alpha$$

### Step 1. Find $\text{Var}(\hat{Z})$

$$\text{Var}(\hat{Z}) = \frac{1}{10-2} \left( 10 \times 3.6695 - 10 \times \frac{33.9667^2}{319.7767} \right) = 0.0769$$

### Step 2. Find critical value of t

$$P(T > t_{n-2, \alpha/2}) = 0.025$$

$$P(T < t_{8, 0.025}) = 0.975$$

$$t_{8, 0.025} = 2.306$$

STUDENT'S t Table

| $\nu$ | 60.0% | 66.7% | 75.0% | 80.0% | 87.5% | 90.0% | 95.0% | 97.5%  | 99.0%  | 99.5%  | 99.9%  |
|-------|-------|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|
| 1     | 0.325 | 0.577 | 1.000 | 1.376 | 2.414 | 3.078 | 6.314 | 12.706 | 31.821 | 63.657 | 318.31 |
| 2     | 0.289 | 0.500 | 0.816 | 1.061 | 1.604 | 1.886 | 2.920 | 4.303  | 6.965  | 9.925  | 22.327 |
| 3     | 0.277 | 0.476 | 0.765 | 0.978 | 1.423 | 1.638 | 2.353 | 3.182  | 4.541  | 5.841  | 10.215 |
| 4     | 0.271 | 0.464 | 0.741 | 0.941 | 1.344 | 1.533 | 2.132 | 2.776  | 3.747  | 4.604  | 7.173  |
| 5     | 0.267 | 0.457 | 0.727 | 0.920 | 1.301 | 1.476 | 2.015 | 2.571  | 3.365  | 4.032  | 5.893  |
| 6     | 0.265 | 0.453 | 0.718 | 0.906 | 1.273 | 1.440 | 1.943 | 2.447  | 3.143  | 3.707  | 5.208  |
| 7     | 0.263 | 0.449 | 0.711 | 0.896 | 1.254 | 1.415 | 1.895 | 2.365  | 2.998  | 3.499  | 4.785  |
| 8     | 0.262 | 0.447 | 0.706 | 0.889 | 1.240 | 1.397 | 1.862 | 2.306  | 2.896  | 3.355  | 4.501  |
| 9     | 0.261 | 0.445 | 0.703 | 0.883 | 1.230 | 1.383 | 1.833 | 2.262  | 2.821  | 3.250  | 4.297  |

### Step 3. Find CI of $\hat{a}$ and $\hat{b}$

$$\sqrt{\frac{\hat{a} - a}{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)}} \sim t(n-2)$$

$$P \left( -t_{n-2,\alpha/2} < \frac{\hat{a} - a}{\sqrt{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)}} < t_{n-2,\alpha/2} \right) = 1 - \alpha$$

$$P \left( \hat{a} - t_{n-2,\alpha/2} \sqrt{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)} < a < \hat{a} + t_{n-2,\alpha/2} \sqrt{\widehat{\text{Var}}(Z) \left( \frac{1}{n} + \frac{\bar{x}^2}{n \text{var}_x} \right)} \right) = 1 - \alpha$$

$$\text{CI of } \hat{a} \text{ at 95\%} = \left( -0.1204 \pm 2.306 \sqrt{0.0769 \left( \frac{1}{10} + \frac{42^2}{10 \times 319.7767} \right)} \right)$$

$$= (-0.6366, 0.3958)$$

$$\sqrt{\frac{\hat{b} - b}{\widehat{\text{Var}}(Z) / n \text{var}_x}} \sim t(n-2)$$

$$P \left( -t_{n-2,\alpha/2} < \frac{\hat{b} - b}{\sqrt{\widehat{\text{Var}}(Z) / n \text{var}_x}} < t_{n-2,\alpha/2} \right) = 1 - \alpha$$

$$P \left( \hat{b} - t_{n-2,\alpha/2} \sqrt{\frac{\widehat{\text{Var}}(Z)}{n \text{var}_x}} < b < \hat{b} + t_{n-2,\alpha/2} \sqrt{\frac{\widehat{\text{Var}}(Z)}{n \text{var}_x}} \right) = 1 - \alpha$$

$$\text{CI of } \hat{b} \text{ at 95\%} = \left( 0.1062 \pm 2.306 \sqrt{\frac{0.0769}{10 \times 319.7767}} \right)$$

$$= (0.0949, 0.1175)$$

d - Test the hypothesis that the slope is zero with significance 5%

### (a) Test on the slope parameter $b$

1. Null Hypothesis  $H_0 : b = b^*$
2. Alternative Hypothesis  $H_1 : b \neq b^*$
3. Test statistic  $T = \frac{\hat{b} - b}{\sqrt{\frac{Var(\hat{Z})}{nVar_x}}}$
4. Distribution of test statistic  $T \sim t(n-2)$
5. Critical Region  $\mathfrak{R}$   $|T| > t_{n-2, \alpha/2}$
6. Evaluate  $T$  under  $H_0$  if  $|t_0| > t_{n-2, \alpha/2} \Rightarrow$  Reject
7.  $p$ -value  $p(|T| > t_0)$

This can easily be transformed to one-sided test where the alternative hypothesis could be  $H_1 : b > b^*$

### Step 1. Define $H_0$ and $H_1$

$$H_0 : b = 0$$

$$H_1 : b \neq 0$$

### Step 2. Define Test Statistic

$$T = \frac{\hat{b} - b}{\sqrt{\frac{Var(\hat{Z})}{nVar_x}}} \quad \leftarrow \text{at this stage, assume } H_0 \text{ is true, } b=0!$$

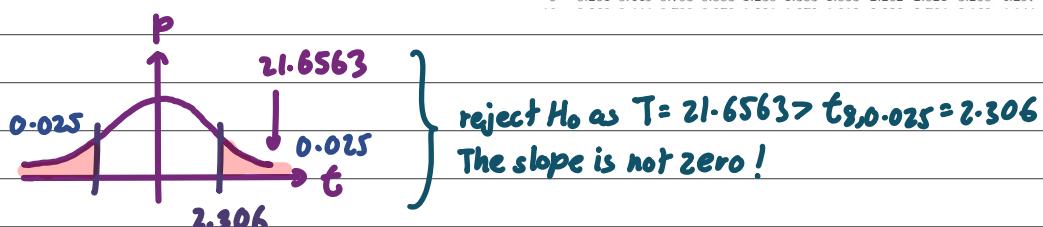
$$= \frac{0.1062}{\sqrt{\frac{0.0769}{10 \times 319.7767}}} \\ = 21.6563$$

### Step 3. Find Critical Region

$$P(T > t_{n-2, \alpha/2}) = 0.025$$

$$P(T < t_{8, 0.025}) = 0.975$$

$$t_{8, 0.025} = 2.306$$



| $\nu$ | 60.0% | 66.7% | 75.0% | 80.0% | 87.5% | 90.0% | 95.0% | 97.5%  | 99.0%  | 99.5%  | 99.9%  |
|-------|-------|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|
| 1     | 0.325 | 0.577 | 1.000 | 1.376 | 2.414 | 3.078 | 6.314 | 12.706 | 31.821 | 63.657 | 318.31 |
| 2     | 0.289 | 0.500 | 0.816 | 1.061 | 1.604 | 1.886 | 2.920 | 4.303  | 6.965  | 9.925  | 22.327 |
| 3     | 0.277 | 0.476 | 0.765 | 0.978 | 1.423 | 1.638 | 2.353 | 3.182  | 4.541  | 5.841  | 10.215 |
| 4     | 0.271 | 0.464 | 0.741 | 0.941 | 1.344 | 1.533 | 2.132 | 2.776  | 3.747  | 4.604  | 7.173  |
| 5     | 0.267 | 0.457 | 0.727 | 0.920 | 1.301 | 1.476 | 2.015 | 2.571  | 3.365  | 4.032  | 5.893  |
| 6     | 0.265 | 0.453 | 0.718 | 0.906 | 1.273 | 1.440 | 1.943 | 2.447  | 3.143  | 3.707  | 5.208  |
| 7     | 0.263 | 0.449 | 0.711 | 0.896 | 1.254 | 1.415 | 1.895 | 2.395  | 2.998  | 3.499  | 4.785  |
| 8     | 0.262 | 0.447 | 0.706 | 0.889 | 1.240 | 1.397 | 1.861 | 2.306  | 2.896  | 3.355  | 4.501  |
| 9     | 0.261 | 0.445 | 0.703 | 0.883 | 1.230 | 1.383 | 1.833 | 2.262  | 2.821  | 3.250  | 4.297  |