

FURTHER Pure MATHEMATICS

4 (PROBABILITY AND STATISTICS)

Amy Khoa
Mdm Lim



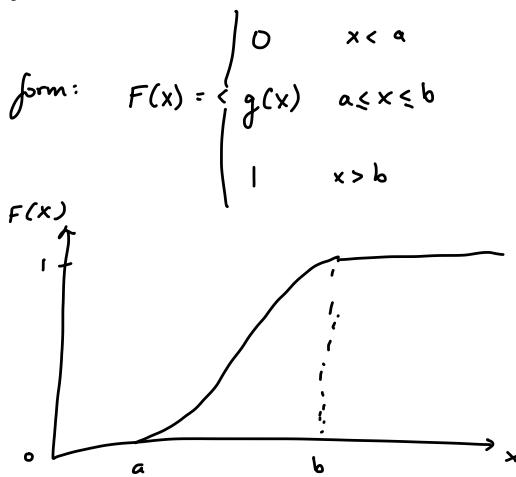
Chapter 1: Distributions

CDF

$$F(x) = P(X < x) = \int_{-\infty}^x f(x) dx.$$

→ shortcut:

in the form: $F(x) = \begin{cases} 0 & x < a \\ g(x) & a \leq x \leq b \\ 1 & x > b \end{cases}$



For a pdf $f(x)$,

$$E(g(x)) = \int_{-\infty}^{\infty} f(x) g(x) dx$$

Transformation of RV

Q: How do you find the paf of $g(y)$ of the transformation $y = g(x)$, w/ a given pdf $f(x)$?

- e.g. $y = x^2$ methods:
 or $y = 2x - 3$ • cdf technique
 or $y = 5 - 2x$ - identify domain of X & Y
 - write $G(Y) = P(Y \leq y)$
 (the CDF of $Y \rightarrow$ the CDF of X)
 - $G'(Y) \rightarrow$ pdf of Y ($g(y)$).
 • pdf technique (not needed).

e.g. 1 (Ex 4 Q2)

$$F(x) = \begin{cases} 0, & x < 0 \\ \frac{1}{10}(x^3 + x), & 0 \leq x \leq 2 \\ 1, & x > 2 \end{cases}$$

Find the CDF of $A = X^3$.

Step 1: find CDF of $f(x)$.

$$\therefore F(x) = \int_{-\infty}^x f(x) dx.$$

Step 2: find $G(a)$, subst $A \rightarrow f(x)$

$$G(a) = P(A < a) = P(X^3 < a)$$

Step 3: make $G(a)$ in the form $P(X < h(a))$, and solve accordingly using $F(x)$

$$= P(X < \sqrt[3]{a})$$

$$G(a) = \frac{1}{10}((\sqrt[3]{a})^3 + \sqrt[3]{a}) = \frac{1}{10}(a + a^{\frac{1}{3}}).$$

Step 4: create the CDF, altering the domain accordingly.
 $\therefore G(a) = \begin{cases} 0, & a < 0 \\ \frac{1}{10}(a + \sqrt[3]{a}), & 0 \leq a \leq 8 \\ 1, & a > 8. \end{cases}$

Step 5: you can find the PDF by differentiating accordingly.

$$g(a) = \begin{cases} \frac{1}{10}(1 + \frac{1}{3}a^{-\frac{2}{3}}), & 0 \leq a \leq 8 \\ 0, & \text{otherwise} \end{cases}$$

e.g. ² $f(x) = \begin{cases} \frac{1}{2}x, & 0 \leq x \leq 2 \\ 0, & \text{otherwise.} \end{cases}$

Find pdf of Y , where $Y = 6X - 3$. \rightarrow calc domain
 $x=0 \rightarrow y=-3$
 $x=2 \rightarrow y=9$.

$$F(x) \rightarrow \int_0^x \frac{1}{2}x \, dx$$

$$= \left[\frac{x^2}{4} \right]_0^x$$

$$= \frac{x^2}{4} - 0$$

$$= \frac{x^2}{4}$$

$$\therefore F(x) = \begin{cases} 0, & x < 0 \\ \frac{x^2}{4}, & 0 \leq x \leq 2 \\ 1, & x > 2. \end{cases}$$

$$G(y) = P(Y < y)$$

$$= P(6X-3 < y)$$

$$= P\left(X < \frac{y+3}{6}\right)$$

$$= F\left(\frac{y+3}{6}\right)$$

$$\therefore G(y) = \begin{cases} 0, & y < -3 \\ \frac{(y+3)^2}{144}, & -3 \leq y \leq 9 \\ 1, & y > 9 \end{cases}$$

$$\therefore g(y) = \begin{cases} \frac{(y+3)^2}{72}, & -3 \leq y \leq 9 \\ 0, & \text{otherwise.} \end{cases}$$

e.g. ³ $f(x) = \begin{cases} 3x^2, & 0 < x < 1 \\ 0, & \text{otherwise.} \end{cases}$

find pdf of Y , $Y = X^2$. $\begin{array}{l} \text{domain:} \\ x=0 \rightarrow y=0 \\ x=1 \rightarrow y=1 \end{array}$

$$F(x) = \int_0^x 3x^2 \, dx$$

$$= \left[x^3 \right]_0^x$$

$$= x^3 - 0 = x^3.$$

$$\therefore F(x) = \begin{cases} 0, & x < 0 \\ x^3, & 0 < x < 1 \\ 1, & x > 1 \end{cases}$$

$$G(y) = P(Y < y)$$

$$= P(X^2 < y)$$

$$= P(X < \sqrt{y})$$

$$= F(\sqrt{y}) = y^{\frac{3}{2}}.$$

$$\therefore G(y) = \begin{cases} 0, & y < 0 \\ y^{\frac{3}{2}}, & 0 < y < 1 \\ 1, & y > 1. \end{cases}$$

$$\therefore g(y) = \begin{cases} \frac{3}{2}y^{\frac{1}{2}}, & 0 < y < 1 \\ 0, & \text{otherwise.} \end{cases}$$

e.g. ⁴ $f(x) = \begin{cases} 1, & 0 < x < 1 \\ 0, & \text{otherwise} \end{cases}$

find pdf of Y . $Y = -\ln x.$ $\begin{array}{l} x=0 \rightarrow y=\infty \\ x=1 \rightarrow y=0 \end{array}$

$$F(x) \rightarrow \int_0^x 1 \, dx$$

$$= [x]_0^x = x.$$

$$\therefore F(x) = \begin{cases} 0, & x \leq 0 \\ x, & 0 < x < 1 \\ 1, & x \geq 1 \end{cases}$$

$$G(y) = P(Y < y)$$

$$= P(-\ln x < y)$$

$$= P(\ln x > -y)$$

$$= 1 - P(\ln x < -y)$$

$$= 1 - P(X < e^{-y}).$$

$$= 1 - F(e^{-y}).$$

$$= 1 - e^{-y}.$$

$$\therefore G(y) = \begin{cases} 0, & y \leq 0 \\ 1 - e^{-y}, & y > 0 \end{cases}$$

$$\therefore g(y) = \begin{cases} e^{-y}, & y > 0 \\ 0, & \text{otherwise.} \end{cases}$$

Chapter 2: Probability Generating Function

Fundamental definitions.

- 1) "generating funct"
→ a funct that stores stats of a prob distn in a compact form.

2) "prob gen func (PGF)"

- used only w/ disc prob distns.
- to obtain one, we introduce a "dummy variable" t , and the PGF will be a func of t .

→ formal defn:

$$G_X(t) = \sum_{\text{all } x} P(X=x) t^x$$

$$G_X(t) = E(t^X)$$

Key results from PGF.

$$\text{① } G_X(1) = 1.$$

$$\text{proof: } G_X(1) = \sum P(X=r) 1^r = \sum P(X=r) = 1.$$

② Power of t is the value of X .
ie $t^k \Leftrightarrow X=k$.

③ If $G_X(t)$ is expanded as a polynomial, the coeff of t^k gives $P(X=k)$.

④ Given the table

x	x_1	x_2	\dots	x_n
$P(X=x)$	p_1	p_2	\dots	p_n

it follows that

$$G_X(t) = \sum_{k=1}^n p_k t^{x_k} = p_1 t^{x_1} + p_2 t^{x_2} + \dots + p_n t^{x_n}.$$

USING A PGF TO FIND μ & σ^2

$$\begin{aligned} \text{⑤ } E(X) &= \sum x_k P(X=x_k) \\ &= x_1 p_1 + x_2 p_2 + \dots + x_n p_n. \\ \cdot E(X^2) &= \sum x_k^2 P(X=x_k) \\ \cdot \text{Var}(X) &= E(X^2) - [E(X)]^2. \end{aligned}$$

We know $G_X(t) = \sum t^x P(X=x)$

$$\therefore G'_X(t) = \sum x t^{x-1} P(X=x) = \sum x^{x-1} x P(X=x)$$

$$t \mapsto 1 \rightarrow G'_X(1) = \sum x P(X=x) = E(X).$$

$$\Rightarrow G'_X(1) = E(X).$$

* formulae in MF19.

PGFs of Special Prob Distns.

① Binomial, $X \sim B(n, p)$ ($n = 1, 2, \dots, n$)

$$\begin{aligned} G_X(t) &= P_0 + P_1 t + P_2 t^2 + P_3 t^3 + \dots + P_n t^n \\ &= {}^n C_0 (1-p)^n + {}^n C_1 (p)(1-p)^{n-1} t + {}^n C_2 (p^2)(1-p)^{n-2} t^2 + \dots + {}^n C_n (p)^n t^n \\ G_X(t) &= (pt + (1-p))^n. \end{aligned}$$

② Poisson, $X \sim Po(\lambda)$ ($n = 0, 1, 2, \dots, \infty$)

$$\begin{aligned} G_X(t) &= P_0 + P_1 t + P_2 t^2 + \dots + P_n t^n \\ &= \frac{e^{-\lambda}}{0!} + \frac{\lambda e^{-\lambda}}{1!} t + \frac{\lambda^2 e^{-\lambda}}{2!} t^2 + \dots + \frac{\lambda^n e^{-\lambda}}{n!} t^n \\ &= e^{-\lambda} \left(1 + \frac{\lambda t}{1!} + \frac{\lambda^2 t^2}{2!} + \frac{\lambda^3 t^3}{3!} + \dots + \frac{\lambda^n t^n}{n!} \right) \xrightarrow{\text{Taylor exp of } e^{\lambda t}} \end{aligned}$$

$$G_X(t) = e^{\lambda(t-1)}$$

③ Geometric, $X \sim Geo(p)$ ($n = 1, 2, \dots, \infty$)

$$\begin{aligned} G_X(t) &= P_1 t^1 + P_2 t^2 + P_3 t^3 + \dots + P_n t^n \\ &= pt + p(1-p)t^2 + p(1-p)^2 t^3 + \dots + p(1-p)^{n-1} t^n \\ &= pt \left(1 + (1-p)t + (1-p)^2 t^2 + \dots \right) \quad (a=1, r=(1-p)t) \\ &= pt \left(\frac{a}{1-r} \right) \end{aligned}$$

$$G_X(t) = \frac{pt}{1-(1-p)t}$$

THE CONVOLUTION THEOREM

$G_{(X+Y)}(t) = G_X(t) G_Y(t)$,

if X & Y are indep.

⇒ the pgf of the Σ of indep rv is = to the product of their pgfs.

Extending this to three or more rv:

$$\Rightarrow G_{\sum X_k}(t) = \prod G_{X_k}(t)$$

LINEAR TRANSFORMATION OF VARIABLES *not stated in syllabus

Given a drv X , & $Y = aX+b$, we can find $G_Y(t)$.

$$\begin{aligned} \Rightarrow G_Y(t) &= E(t^Y) \\ &= E(t^{aX+b}) \\ &= t^b E(t^{aX}) \\ &= t^b E((t^a)^X) \\ &= t^b G_X(t^a) \end{aligned}$$

result: $G_Y(t) = t^b G_X(t^a). \quad (Y=aX+b)$

Chapter 3:

Sampling, Confidence Interval and Hypothesis Testing

INTRODUCTION

- A "population" includes all the elements from a set of data.
- A "sample" consists of ≥ 1 observations from the population.
- "Sampling" involves extrapolating data from a sample of the population to the whole population.
- A "parameter" is a measurable characteristic of a population.
- A "statistic" is a measurable characteristic of a sample.

SAMPLING DIST^N OF SAMPLE MEAN

We use " \bar{X} " to denote the mean of a sample of a population.

\Rightarrow this is a drv, with its possible values being the mean of the different possible samples of the population.

① Mean of Sample Mean

The mean of the sampling distⁿ of the mean will be equivalent to the mean of the population.

$$\text{i.e. } E(\bar{X}) = E(X) = \mu.$$

② Variance of Sample Mean

The variance of the sampling distⁿ of the mean is equivalent to the variance of the population divided by the sample size.

$$\text{i.e. } \text{Var}(\bar{X}) = \frac{\text{Var}(X)}{n} = \frac{\sigma^2}{n}.$$

Central limit theorem (CLT)

The central limit theorem states that, given $n \geq 30$, any sampling distⁿ tends towards $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$.

e.g. for a population x_1, x_2, x_3, x_4 , and sample size 2,
 $\bar{X} \in \left\{ \frac{x_1+x_2}{2}, \frac{x_1+x_3}{2}, \frac{x_1+x_4}{2}, \dots, \frac{x_3+x_4}{2} \right\}$.

• Sampling Distⁿs of some special prob distⁿs

① Normal, $X \sim N(\mu, \sigma^2)$
 $\Rightarrow \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

*don't forget the $(\frac{1}{n})$!!

② Poisson, $X \sim Po(\lambda)$

$$\Rightarrow \bar{X} \approx N\left(\lambda, \frac{\lambda}{n}\right)$$

③ Geometric, $X \sim Geo(p)$

$$\Rightarrow \bar{X} \approx N\left(\frac{1}{p}, \frac{q^2}{pn}\right)$$

④ Binomial, $X \sim B(n, p)$

$$\Rightarrow \bar{X} \approx N\left(np, \frac{np(1-p)}{n}\right)$$

*these are approximations as we are estimating drv to a crv.

ESTIMATION

"Estimation" refers to the process in statistics where we can deduce characteristics about a population from measuring data from a sample.

\Rightarrow this is done via the use of "sample statistics".

ESTIMATION METHODS

① Point estimate

\Rightarrow a single value of a statistic.

\Rightarrow eg \bar{x} (sample mean) $\approx \mu$ (pop mean)

② Interval estimate

\Rightarrow defined by 2 numbers.

\Rightarrow eg $a < \mu < b$.

POINT ESTIMATION

Conceptualisation:
Let our sample contains n elements, x_1, x_2, \dots, x_n , where $x_i \sim N(\mu, \sigma^2)$ for $1 \leq i \leq n$.

\Rightarrow denote \bar{x} as the sample mean, and

\Rightarrow denote $\hat{\sigma}^2$ as the sample variance.

Keep in mind \bar{x} and $\hat{\sigma}^2$ are also random variables, but are influenced by the values of x_i .

\Rightarrow Hence, by defn.

$$\bar{x} = \frac{\sum x_i}{n} \quad \text{and} \quad \hat{\sigma}^2 = \frac{\sum x_i^2}{n} - \left[\frac{\sum x_i}{n} \right]^2.$$

The unbiased estimate of the population mean

\Rightarrow The value of $E(\bar{x})$, or $\hat{\mu}$, give us our estimate for the population mean.
 \Rightarrow from a sample we are making inferences about a population.

$$\Rightarrow E(\bar{x}) = E\left(\frac{\sum x_i}{n}\right)$$

$$= \frac{1}{n} E(\sum x_i)$$

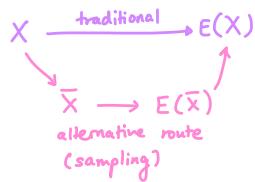
$$= \frac{1}{n} E(x_1 + x_2 + \dots + x_n)$$

$$= \frac{1}{n} n E(x)$$

$$= E(x)$$

$$= \mu.$$

Conclusion := $\hat{\mu} = E(\bar{x}) = \mu$
 \Rightarrow mean of sample mean is equal to mean of population.



The unbiased estimate of the population variance

\Rightarrow Similarly, $E(\hat{\sigma}^2)$ will give us an estimate for the population variance.

$$\begin{aligned} \Rightarrow E(\hat{\sigma}^2) &= E\left(\frac{\sum x_i^2}{n} - \left[\frac{\sum x_i}{n}\right]^2\right) \\ &= \frac{1}{n} E(\sum x_i^2) - E(\bar{x}^2) \\ &= \frac{1}{n} E(\sum x_i^2) - \left(\frac{\sigma^2}{n} + \mu^2\right) \quad \text{recall } \text{Var}(\bar{x}) = \frac{\sigma^2}{n} \\ &= \frac{1}{n} E(x_1^2 + x_2^2 + \dots + x_n^2) - \left(\frac{\sigma^2}{n} + \mu^2\right) \quad = E(\bar{x}^2) - \left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= \frac{1}{n} E(n\bar{x}^2) - \left(\frac{\sigma^2}{n} + \mu^2\right) \quad \text{recall } \text{Var}(X) = E(X^2) - [E(X)]^2 \\ &= E(\bar{x}^2) - \left(\frac{\sigma^2}{n} + \mu^2\right) \quad \Rightarrow E(\bar{x}^2) = \frac{\sigma^2}{n} + \mu^2. \\ &= (\sigma^2 + \mu^2) - \left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= \frac{1}{n} (n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2) \\ &= \frac{n-1}{n} \sigma^2. \quad (\approx \sigma^2) \end{aligned}$$

Conclusion := $E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2$.
 \star the sample variance is not equal to the population variance!

\Rightarrow How do we fix this?

$$\text{Notice } \frac{n}{n-1} E(\hat{\sigma}^2) = \sigma^2$$

$$E\left(\frac{n}{n-1} \hat{\sigma}^2\right) = \sigma^2.$$

\Rightarrow hence, the unbiased estimate of the population variance is actually $\frac{n}{n-1} \hat{\sigma}^2$, which is also denoted as $\underline{s^2}$.

$$\text{ie } s^2 = \frac{n}{n-1} \hat{\sigma}^2.$$

INTERVAL ESTIMATION (FOR μ ONLY)

Q: Interval estimates are preferred over point estimates, as it also provides a "confidence level" for the estimate.

⇒ these are called "confidence intervals". eg 90% confidence interval ⇔ pop mean falls in interval.

Relationship bw pt est, confidence interval

\bar{z}



From the diagram, it is apparent that:

$$\begin{aligned} \text{① } P(\bar{x}_1 < \bar{x} < \bar{x}_2) &= 1 - \alpha. \\ &\Rightarrow P(\bar{x} > \bar{x}_2) = P(\bar{x} < \bar{x}_1) = \frac{\alpha}{2}. \\ \text{② } \bar{x}_2 &= -\bar{x}_1. \end{aligned}$$

Hence:

$$P\left(\frac{\bar{z}}{\hat{\sigma}} > \frac{\bar{x}_2 - \bar{x}}{\hat{\sigma}}\right) = P\left(\frac{\bar{z}}{\hat{\sigma}} < \frac{\bar{x}_1 - \bar{x}}{\hat{\sigma}}\right) = \frac{\alpha}{2}.$$

$$\hat{\sigma} = \frac{\sigma}{\sqrt{n}}$$

$$\Rightarrow P\left(\frac{\bar{z}}{\frac{\sigma}{\sqrt{n}}} > \frac{\bar{x}_2 - \bar{x}}{\frac{\sigma}{\sqrt{n}}}\right) = P\left(\frac{\bar{z}}{\frac{\sigma}{\sqrt{n}}} < \frac{\bar{x}_1 - \bar{x}}{\frac{\sigma}{\sqrt{n}}}\right) = \frac{\alpha}{2}.$$

$$\text{Let } \frac{\bar{x}_2 - \bar{x}}{\frac{\sigma}{\sqrt{n}}} = z \Rightarrow \frac{\bar{x}_1 - \bar{x}}{\frac{\sigma}{\sqrt{n}}} = -z. \quad \text{i.e.}$$

$$\Rightarrow \bar{x}_2 = \bar{x} + z \frac{\sigma}{\sqrt{n}}, \quad \bar{x}_1 = \bar{x} - z \frac{\sigma}{\sqrt{n}}.$$

Conclusion: The confidence interval is

$$\boxed{\bar{x} - z \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z \frac{\sigma}{\sqrt{n}}}$$

alternative ways to represent:
1) $(\mu - z \frac{\sigma}{\sqrt{n}}, \mu + z \frac{\sigma}{\sqrt{n}})$
2) $\mu \pm z \frac{\sigma}{\sqrt{n}}$

$$\text{(where } z = \phi^{-1}\left(\frac{\alpha+100}{200}\right) \text{ or } \alpha = 90\%, \therefore z = \phi^{-1}(0.95))$$

⇒ or:
the true value of the population mean, μ , lies in this confidence interval.

* remember:

μ = pop mean

\bar{x} = sample mean

$\hat{\mu}$ = mean of the sample mean

dist^n

Cases

① μ & σ known

For $n \leq 30$, we make 2 assumptions:

- distⁿ is normally distributed
- pop. variance is known

For $n \geq 30$, we apply CLT.

The confidence interval is just the direct application of the formula:

$$\bar{x} - z \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z \frac{\sigma}{\sqrt{n}}$$

② μ is known, σ is unknown

We must approximate σ with s if we are not given σ directly, given that $n \leq 30$.

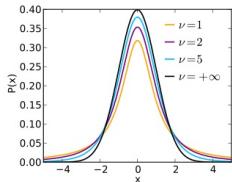
⇒ however, this introduces more uncertainty!

⇒ to compensate for this, instead of using a z -distribution, we use a t -distⁿ instead for our calculations.

t_{n-1} = critical value of t distⁿ w/ $\nu = n-1$ & area of $\frac{\alpha}{2}$ in each tail.

t-DISTNS VS Z-DISTNS

Q: We use t-distns to estimate a confidence interval when $n \leq 30$, σ is unknown, & $X \sim N(\mu, \sigma^2)$.



Characteristics

① Bell shaped and symmetrical but with fatter tails than the normal.

② Degree of freedom

- denoted by ν
- look at table in MF19
- when $\nu \rightarrow \infty$, t distⁿ → normal.

① If t distⁿ is used:

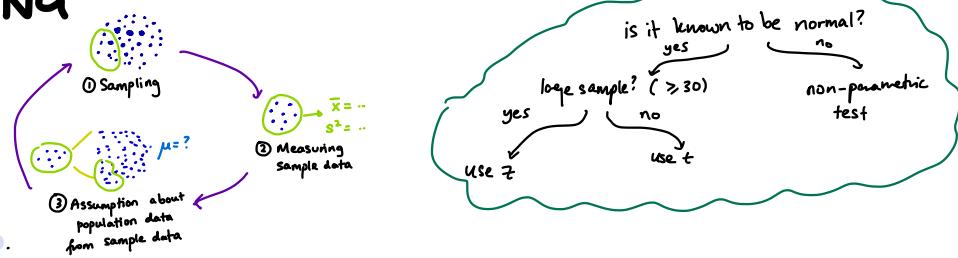
$$z = \phi_{\nu=?}^{-1}\left(\frac{\alpha+100}{200}\right).$$

② If two/paired sample is used:

$$\begin{aligned} \frac{\sigma}{\sqrt{n}} &= s = \sqrt{s_c^2} \text{ or } \sqrt{s_p^2} \\ &= \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad = \sqrt{\frac{n_1 \bar{s}_1^2 + n_2 \bar{s}_2^2}{n_1 + n_2 - 2}} \end{aligned}$$

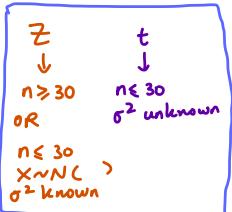
HYPOTHESIS TESTING

By examining characteristics (mean / variance) of a sample, we can make inferences about the population the sample was derived from. ↳ hypothesis testing can help us evaluate whether there is sufficient evidence to justify a given assumption about a population.



ONE SAMPLE Z & t-TEST

An one sample test is used to determine whether a population parameter is significantly different than some hypothesised value.



- 1 Steps
 - 1 State variable, check assumptions
 - 2 State the null (H_0) & alternative (H_1) hypothesis
 - 3 State sampling dist according to H_0 .
 - 4 State level of test.
 - 5 Find the critical value of the test, via the decision rule.
 - 6 Perform the necessary calculations.

$$\begin{cases} Z_t \text{ & } Z_c \text{ comparison} \\ Z_t \text{ & } \bar{X}_c \text{ comparison} \\ P_{\alpha} \text{ & } \alpha \text{ comparison.} \end{cases}$$
 \Rightarrow write down the rejection criteria.
 - 7 Form a conclusion.

If H_0 rejected:
 "we reject H_0 at (value of α) level of significance. There is sufficient evidence to conclude that (statement of the alternative hypothesis).
 (and " H_0 accepted" statement follows same structure).

Worked example (Z)

Example 6 Boys of a certain age are normally distributed and known to have a mean weight of $\mu = 85$ pounds. A complaint is made that the boys living in a municipal children's home are underfed. As one bit of evidence, $n = 25$ boys (of the same age) are weighed and found to have a mean weight of 80.94 pounds. It is known that the population standard deviation σ is 11.6 pounds. Test at 5% significant level, what should be concluded concerning the complaint?

Step 1 (State variable of population)

Let X be weight of boys/ pounds.

$$\therefore X \sim N(\mu, 11.6^2)$$

Step 2 (State hypotheses)

$$H_0 : \mu = 85$$

$$H_1 : \mu < 85$$

Step 3 (State test & variable of sample, assuming H_0 is true)

$\rightarrow Z$ -test

$$\therefore \bar{X} \sim N(85, \frac{11.6^2}{25})$$

Step 4 (State significance level)

$$\alpha = 0.05.$$

Step 5 & 6 (comparison & conditions for rejection) (choose 1 method)

i) Compare Z_t to Z_c (critical value)

$$Z_c = \Phi^{-1}(1-0.05)$$

$$= \Phi^{-1}(0.95)$$

$$= -1.645.$$

* We reject H_0 IF

$$Z_t \leq -1.645.$$

$$\begin{aligned} Z_t &= \frac{\bar{X} - \hat{\mu}}{\hat{\sigma}} && (\bar{X} = \text{test statistic of mean}) \\ &= \frac{80.94 - 85}{\frac{11.6}{\sqrt{25}}} && (\hat{\mu} = \text{expected mean}) \\ &= -1.75 && (-1.75 < -1.645) \Rightarrow \text{can reject } H_0! \end{aligned}$$

Step 7 (State Conclusion)

"Hence, H_0 is rejected at 5% significance level.
 It subsequently follows that there is sufficient evidence to conclude the complaint that the boys living in the municipal children's home are underfed is valid."

Worked example (T)

Example 7 Monarch butterflies are bred at a butterfly farm. Monarch butterflies should grow to have a mean wingspan of 9.4 cm. The breeders are concerned that their butterflies are growing as well as they could be, so a sample of six Monarch butterflies is taken and their wingspans (in cm), These are recorded as: 8.8, 9.6, 9.2, 9.1, 9.9, 8.7. Assuming that the wingspans are normally distributed and using a 10% significance level, investigate whether the wingspans of the butterflies are less than 9.4 cm.

Step 1 (variable)

(let X be wingspan of butterflies / cm)

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

Step 2 (hypothesis)

$$H_0 \rightarrow \mu = 9.4$$

$$H_1 \rightarrow \mu < 9.4$$

Step 3 (given H_0 true, find \bar{X})

$\rightarrow t$ -Test $\because \sigma^2$ unknown

$$\bar{X} = \frac{\sum x}{n} = \frac{55.3}{6} = 9.21666$$

$$\sigma^2 = \frac{\sum x^2 - \bar{x}^2}{n} = \frac{510.75}{6} - 9.2166^2 = 0.178055$$

$$= 0.178055.$$

Step 5 & 6 (comparison)

$$P(T \leq t_c) = 0.1, v = 5$$

$$\Rightarrow P(T \leq -t_c) = 0.9$$

$$\therefore -t_c = 1.476$$

$$t_c = -1.476. \text{ (crit value).}$$

\Rightarrow We reject H_0 if $t_t < -1.476$.

$$t_t = \frac{\bar{X} - \hat{\mu}}{\left(\frac{\sigma}{\sqrt{n}}\right)} \quad (\bar{X} = \text{mean obtained from sample})$$

$$= \frac{9.21666 - 9.4}{\sqrt{0.178055}} = -0.1715 (> -1.476).$$

Step 7 (conclusion)

From the results, we accept H_0 using a 10% significance level.

There is hence insufficient data to suggest the mean wingspan of the Monarch butterflies is less than 9.4cm.

Step 4 (state or)

$$\rightarrow \alpha = 0.1$$

TWO SAMPLE TEST

We perform a two-sample test on two random samples, each from 2 different populations.
For our syllabus, we test the difference between the 2 means.

Variance of the two samples

There are 2 ways in obtaining the variance of the two samples:

① Combined variance

case 1 variances are known

$$\Rightarrow \sigma_{\text{comb}}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

σ_1^2 = variance of pop. 1
 σ_2^2 = variance of pop. 2
 n_1 = sample size of 1
 n_2 = sample size of 2.

case 2 variances are unknown

$$\Rightarrow s_{\text{comb}}^2 \approx \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}$$

s is replaced by the Sp. (pooled variance).

② Pooled Variance

We can calculate a pooled variance, otherwise known as the two-sample estimate of a common variance, when the means of the populations are different, but the variances can be assumed to be the same.

$$(MF19) \quad s^2 = \frac{\sum(x_1 - \bar{x}_1)^2 + \sum(x_2 - \bar{x}_2)^2}{n_1 + n_2 - 2}$$

Alternate forms

$$\Rightarrow s^2 = \frac{n_1 \hat{\sigma}_1^2 + n_2 \hat{\sigma}_2^2}{n_1 + n_2 - 2}$$

$$\left(\frac{\sum(x - \bar{x})^2}{n} = \hat{\sigma}^2 \right) \Rightarrow s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

DIFFERENCE BW TWO MEANS

We can use a two-sample test to test the significance of a difference between two population means.

Steps

① Define the variables.

$$X_1 \sim N(\mu_1, \sigma^2) \quad X_2 \sim N(\mu_2, \sigma^2)$$

② State the null and alternative hypotheses.

(One-tailed) $H_0: \mu_1 = \mu_2$ (two-tailed) $H_0: \mu_1 = \mu_2$

$H_1: \mu_1 > \mu_2$ — case 1 $H_1: \mu_1 \neq \mu_2$ — case 3
 or $\mu_1 < \mu_2$ — case 2

③ Compute the difference in the sample means.

$$d = \bar{X}_1 - \bar{X}_2$$

$$= \frac{\sum X_1}{n_1} - \frac{\sum X_2}{n_2}$$

④ Compute the combined variance.

i) If σ^2 is known:
 \Rightarrow compute directly.

$$\sigma_c^2 = \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}$$

$$= \frac{(n_1 + n_2)}{n_1 n_2} \sigma^2$$

* rmb variances of 1 & 2 are the same!

ii) If σ^2 is unknown:
 \Rightarrow compute pooled variance first,
 \Rightarrow then compute the combined variance. *estimate.

$$S_p^2 = \frac{n_1 \hat{\sigma}_1^2 + n_2 \hat{\sigma}_2^2}{n_1 + n_2 - 2}$$

$$S_c^2 = \frac{S_p^2}{n_1} + \frac{S_p^2}{n_2}$$

$$= \frac{(n_1 + n_2)}{n_1 n_2} S_p^2$$

⑤ Using the given significance level, α , calculate critical values.

i) If σ^2 is known, use the z-distn.

$$z_c = \phi^{-1}(1-\alpha) \quad (\#1)$$

$$= -\phi^{-1}(1-\alpha) \quad (\#2)$$

$$= \pm \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \quad (\#3)$$

ii) If σ^2 is unknown, use the t-distn, with $v = n_1 + n_2 - 2$.

$$t_c = \phi_{v=n_1+n_2-2}^{-1}(1-\alpha) \quad (\#1)$$

$$= -\phi_{v=n_1+n_2-2}^{-1}(1-\alpha) \quad (\#2)$$

$$= \pm \phi_{v=n_1+n_2-2}^{-1}\left(1 - \frac{\alpha}{2}\right) \quad (\#3)$$

⑥ Using values from the samples, calculate the test values. (Assume H_0 is correct $\Rightarrow \mu_1 = \mu_2 = 0$)

i) If σ^2 is known

$$z_t = \frac{d - (0)}{\sigma_c}$$

$$z_t = \frac{d}{\sigma_c}$$

ii) If σ^2 is unknown

$$t_t = \frac{d - (0)}{s_c}$$

$$t_t = \frac{d}{s_c}$$

⑦ Check to see whether z_t / t_t lies within the rejection region determined by z_c / t_c .

You can then subsequently write your conclusion.

(lies in region) From the results, we can reject H_0 with a (100) $\%$ significance level.

It hence follows that there is sufficient evidence to suggest ..

(lies out of region) From the results, we can accept H_0 with a (100) $\%$ significance level.

It hence follows that there is insufficient evidence to suggest ..

Key Assumptions

- $X_1 \sim NC$, $X_2 \sim NC$
- Values from the samples are independent from each other
- $\sigma_1^2 = \sigma_2^2$

PAIRED SAMPLE TEST

- We can use a paired sample test when the two samples are correlated; for example, before & after, or when the samples occur in pairs.
- $\star z\text{-test: } n \geq 30$
- $\star t\text{-test: } n < 30$
- objective: conclude the significance of the difference bw the two means.

Key assumptions

- For paired sample t-test, assume Z-dist!.
- \star Variance of the two samples are the same.
- Treat the "difference" as a one-sample test.
(so like $\bar{X} \sim N(\hat{\mu}, \frac{\sigma^2}{n})$)

DIFFERENCE BW TWO MEANS IN PAIRED SAMPLE TEST.

- State variables, hypotheses & the significance level.
 $H_0: \mu_1 - \mu_2 = c$
- Tabulate the difference in the sample test means. $(\hat{\mu}_1 - \hat{\mu}_2)$.
- Find t_c or z_c .

$$t_c = \Phi^{-1}_{v=n_1+n_2-2} (\text{sig lvl}) \quad z_c = \Phi^{-1} (\text{sig lvl})$$

- Find t_t or z_t .

$$t_t = \frac{(\hat{\mu}_1 - \hat{\mu}_2) - c}{\left(\frac{s_d}{\sqrt{n}}\right)}$$

- Compare and make your conclusions.

$$z_t = \frac{(\hat{\mu}_1 - \hat{\mu}_2) - c}{\left(\frac{s_d}{\sqrt{n}}\right)}$$

* s_d = unbiased est. of variance of the diff of 2 samples.

σ^2 same
 $s_d = s_{\text{pool}}$

σ^2 different
 $s_d = s_{\text{comb}}$

Chapter 4: Chi-Square Tests

Chi-squared (χ^2) tests are primarily used for:

- ① to test for the "goodness of fit" of an observed distⁿ to a theoretical one.
- ② to test the null hypothesis that the variables are independent.

χ^2 STATISTIC

A χ^2 statistic is a measurement of how expectations compare to results.

* data used to calculate this

MUST be:

- 1) random;
- 2) raw;
- 3) mutually exclusive;
- 4) drawn from independent variables; and
- 5) drawn from a large enough sample.

χ^2 GOODNESS OF FIT TEST

We use the χ^2 goodness of fit test when:

- sampling is purely random;
- the variable is categorical; → can only take discrete values.
- the expected value of the number of sample observations in each level of the variable is ≥ 5 . → the different values the variable can take.

STEPS IN χ^2 GOF TEST.

① State the hypotheses.

H_0 : data is consistent w/ a specified distⁿ.

H_1 : data is not consistent w/ a specific distⁿ.

② Formulate an analysis plan.

The analysis plan describes how to use sample data to accept or reject H_0 .

a) State the significance level.

→ any value $\in (0, 1)$

→ usually 0.01, 0.05, 0.10.

b) Use the χ^2 GDF test.

→ does the observed sample frequencies differ significantly from the expected frequencies?

③ Analyse sample data.

→ expected freq of a value, $E_i = np_i$

(n = total freq, p_i = prob of value)

→ χ^2 test statistic, $\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$

→ degrees of freedom, $\nu = k-1 - \# \text{ of parameters est.}$

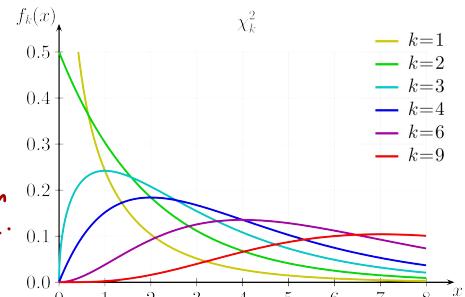
(k = # of levels of categorical variable)

CHI-SQUARE DIST^N

The χ^2 -distⁿ is the distⁿ of the sum of squared std. normal deviates.

★ std normal deviate = a random sample from the std. normal distⁿ.

- degrees of freedom, k , is equal to the number of std. normal deviates being summed.
- ⇒ this is also equal to the "level" of the variable.



Worked example 1

Acme Toy Company prints baseball cards. The company claims that 30% of the cards are rookies, 60% veterans but not All-Stars, and 10% are veteran All-Stars.

Suppose a random sample of 100 cards has 50 rookies, 45 veterans, and 5 All-Stars. Is this consistent with Acme's claim? Use a 0.05 level of significance.

① State hypotheses.

H_0 : consistent

H_1 : inconsistent.

② Formulate analysis plan.

$\alpha = 0.05$.

χ^2 GDF Test.

③ Analyse sample data.

	O_i	$E_i = np_i$	$\frac{(O_i - E_i)^2}{E_i}$
rookies	50	$100 \times 0.3 = 30$	13.33
veterans	45	$100 \times 0.6 = 60$	3.75
all-stars	5	$100 \times 0.1 = 10$	2.50
	$\Sigma 100$	100	19.58

★ if $E_i < 5$, combine the row w/ the row below.
Repeat this until $E_i \geq 5$.

eg E_i $\begin{matrix} 0.12 \\ 0.34 \\ 4.6 \end{matrix} \rightarrow \begin{matrix} 0.46 \\ 4.6 \end{matrix} \Rightarrow \begin{matrix} 5.06 \\ 4.6 \end{matrix}$

$$\therefore \chi^2_t = 19.58$$

④ Calculate critical value of χ^2 .

$$\nu = 2 \quad p = 0.95$$

(from the table) $\therefore \chi^2_c = 5.991$. ($< \chi^2_t$)

⑤ Conclusion.

$$\chi^2_t > \chi^2_c$$

⇒ reject H_0

⇒ so proportion of cards proposed in the Q is wrong.

④ Calculate the critical value of χ^2 .

★ use the table in data booklet.

⑤ Make your conclusions.

★ Compare χ^2_t and χ^2_c .

If $\chi^2_t > \chi^2_c \Rightarrow$ reject H_0
 \hookrightarrow data not consistent w/ distⁿ.

If $\chi^2_t < \chi^2_c \Rightarrow$ accept H_0
 \hookrightarrow data consistent w/ distⁿ.

TEST FOR VARIABLE INDEPENDENCE USING A CONTIGENCY TABLE

In statistics, we can use a **contingency table** (also called a cross tabulation) as a way of displaying how one categorical variable is "distributed" across another categorical variable. (see right for example).

→ using χ^2 tests, we can determine whether any relationship is present between the variables (ie test for independence).

METHOD

Q: Given the following contingency table, determine whether a given person's income level affects their method of transport, using a 5% significance level.

		Method of Transport			Total
		Car	Public	Self	
Income Level	Small	58	21	36	115
	Average	199	49	64	312
Total	462	102	129	693	

① State hypotheses.

H_0 : income level & method of transport are independent

H_1 : ... " " " not independent.

② For each entry in the contingency table, calculate the expected frequency.

∴ Since we assume H_0 to be true,

for a given entry "i" with row attr "A" & column attr "B", and for a d.r.v X , that takes a random entry in the table;

$$\text{expected freq. } E_i = P(X=A) \times P(X=B) \times \Sigma f \\ = \frac{\text{row total}}{\text{grand total}} \times \frac{\text{column total}}{\text{grand total}} \times \text{grand total}$$

$$\Rightarrow E_i = \frac{\text{row total} \times \text{column total}}{\text{grand total}}$$

E_i	Method of transport		
	Car	Public	Self
Small	76.666	16.926	21.407
Avg	208.000	45.922	58.077
Large	177.333	39.152	49.515

Example for further understanding

A	a_1	a_2	a_3	total
B	100	200	100	400
b_1	350	400	250	1000
total	450	600	350	1400

The table shows how variable "A" is "distributed" across variable "B".

③ For each entry "i", work out the squared std. normal deviate ($= \frac{(O_i - E_i)^2}{E_i}$).

income level	method of transport	O_i	E_i	$\frac{(O_i - E_i)^2}{E_i}$
Small	car	58	76.666	5.5707
	public	21	16.926	0.9806
	self	36	21.407	9.9479
average	car	199	208.000	0.38942
	public	49	45.922	0.20630
	self	64	58.077	0.60406
large	car	205	177.333	4.3165
	public	32	39.152	1.3065
	self	29	49.515	8.4997

★ Same idea;
if $E_i < 5$,
combine the rows.

$$\therefore \chi^2_t = 31.827.$$

④ Analyse the data;

work out χ^2_c and χ^2_t , and their relationship. (< or >)

∴ Degree of freedom = (# of rows - 1)(# of col. - 1)

$$(from ③) \chi^2_t = 31.827.$$

$$\alpha = 0.95 \\ \nu = (3-1)^2 = 4 \\ \therefore \chi^2_c = 15.51.$$

$$\Rightarrow \chi^2_t > \chi^2_c.$$

⑤ Form appropriate conclusions.

⇒ Hence, we can reject H_0 with a 5% significance level.

⇒ Consequently, the method of transport and the income level are not independent.

Chapter 5:

Non-Parametric Tests

Compared to parametric tests, which assume the data it is applied upon has a normal distⁿ, non-parametric tests do not assume anything about the data's underlying distⁿ.

→ usually, we use these tests when we know the population data does not have a normal distⁿ.

* Hence, a non-parametric test is any hypothesis test that do not rely on (any assumptions about) the shape of the distⁿ, or any population parameters (e.g. μ , σ^2 etc)

Why do we use non-parametric tests?

• We generally use non-parametric tests when we know $X \sim N(\mu, \sigma^2)$; however, there are other instances where we use them over their parametric counterparts:

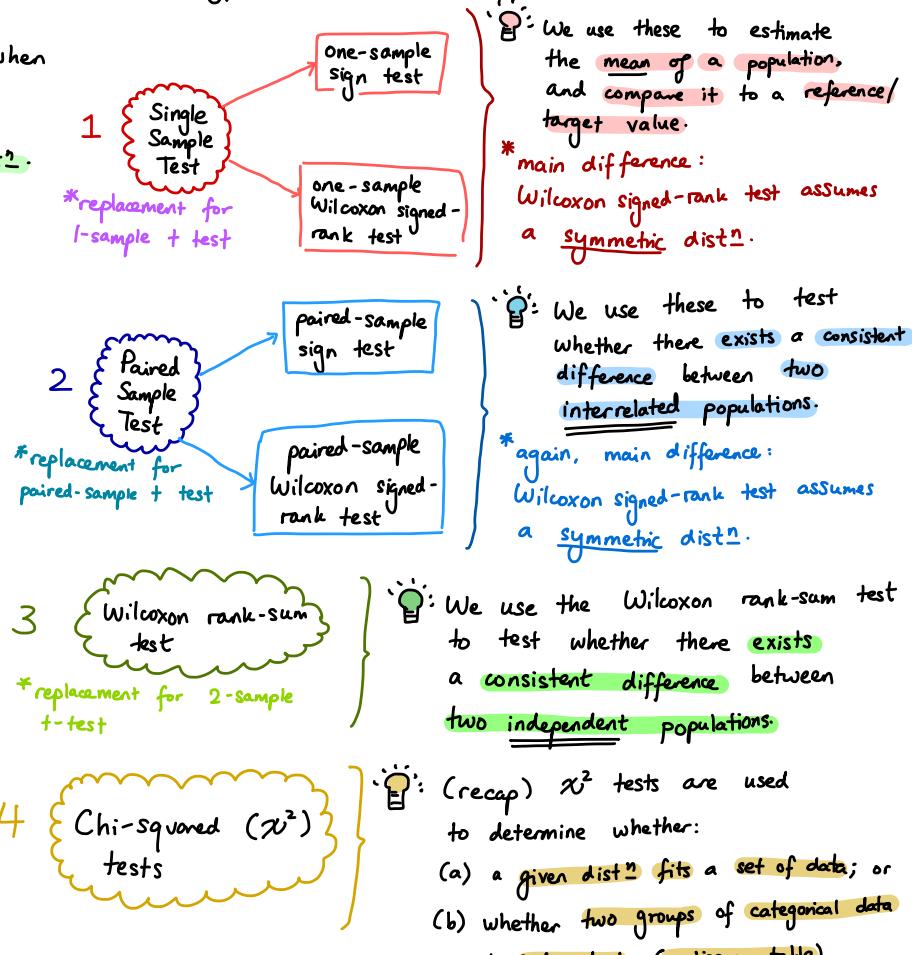
- ① The sample size is lll;
- ② The data contains irremovable outliers;
- ③ You want to test for the median rather than the mean;

Note: we generally do this for extremely skewed distⁿs.

Why? ⇒ mean is affected by outliers.

TYPES OF NON-PARAMETRIC TESTS

• Generally, non-parametric tests can be classified into several types:



• We use these to estimate the mean of a population, and compare it to a reference/target value.
* main difference: Wilcoxon signed-rank test assumes a symmetric distⁿ.

• We use these to test whether there exists a consistent difference between two interrelated populations.
* again, main difference: Wilcoxon signed-rank test assumes a symmetric distⁿ.

• We use the Wilcoxon rank-sum test to test whether there exists a consistent difference between two independent populations.

• (recap) χ^2 tests are used to determine whether:
(a) a given distⁿ fits a set of data; or
(b) whether two groups of categorical data is independent (contingency table).

Pros and cons of non-parametric tests over parametric tests

Advantages

- ✓ fewer assumptions
- ✓ caters for lll sample sizes
- ✓ caters for more data types

Disadvantages

- ✗ less powerful than parametric tests (if assumptions are inviolate)
- ✗ more tedious to calculate by hand

ONE SAMPLE TESTS

We use one sample non-parametric tests to compare the test median of a set of data against a hypothesised median value.

→ two types : { sign test; and Wilcoxon signed-rank test.

1-SAMPLE SIGN TEST

Requirements:

- 1) each pair of data pts is ordered; ie for $\{u_0, u_1, \dots, u_k\}$, $u_i > u_{i-1} \quad \forall 1 \leq i \leq k$

* Example: Let a sample obtained from a population be $\{1, 2, 2, 4, 5, 7, 10, 12, 15, 17\}$. Use a non-parametric test to test whether the median = 8, at a 5% significance level.

Method (One-tailed)

① State the hypotheses.

$$H_0: m = M_0$$

$$H_1: \begin{cases} m > M_0 \\ \text{case #1} \end{cases} \quad \text{or} \quad \begin{cases} m < M_0 \\ \text{case #2} \end{cases}$$

② (Case #1) Find S^+ ; (Case #2) Find S^- . } = S (test statistic). * discard any elements = to the median.

$$(O_{i_1} \text{ or } O_{i_2} = S).$$

③ If we consider H_0 to be true,

$$P(x \leq m_0) = P(x \geq m_0) = \frac{1}{2}. \quad (\text{defn of the median.})$$

Hence, if we let

$$X \mapsto \# \text{ of observations (case #1)} > m_0 \\ (\text{case #2}) < m_0, \\ \text{out of } n,$$

we see $X \sim B(n, \frac{1}{2})$, as all observations are independent, and the chance any given observation will be larger than the median is equal to the chance the observation will be smaller than the median.

Given a significance level α , we say that we reject H_0 if if $P(S^+ < s_t) < \alpha$ (case #1) or $P(S^- < s_t) < \alpha$ (case #2).

Note | #1 For 2-tailed tests,

we check

$$2P(X \geq S) \leq \alpha,$$

and if this is true we reject H_0 .

| #2 For large samples ($n \geq 30$), we can use a normal approximation;

$$X \approx N(np, np(1-p)).$$

Hence, $P(X \geq s_t \leq S)$

$$= P(Z \geq s_t - \frac{S - \frac{n}{2}}{\sqrt{\frac{1}{4}n}}).$$

1-SAMPLE WILCOXON SIGNED-RANK TEST

Example

Wilcoxon signed rank test
Example 4
The weights (in kg) of ten randomly selected Spanish mackerel are recorded:
1.6 1.1 2.1 2.4 2.2 2.9 2.6 2.3 2.7 1.9
Test, at the 5% significant level, whether the median weight is greater than 1.8 kg.

Assumptions:

1) symmetric dist \Rightarrow mean = median.

2) continuous dist

3) data is independent.

Method

① State hypotheses.

$$H_0: \text{pop. med} = M$$

$$H_1: \text{pop. med} \begin{cases} \neq M \\ > M \\ < M. \end{cases}$$

② Calculate the test statistic:

a) Calculate the absolute differences bw each sample value & the hypothesised median.

$$\text{ie } |x_1 - M|, |x_2 - M|, \dots, |x_n - M|.$$

b) Rank each values from 1 to n , giving the lowest rank to the smallest absolute difference.

$$\therefore x_t = \min(W^+, W^-) \\ = \min(46, 9) \\ = 9.$$

c) Calculate $\underline{W}^+ = \text{the sum of ranks of the sample values} > M$;

& $\underline{W}^- = \text{the sum of ranks of the sample values} < M$.

d) The test (signed rank) statistic, $x_t = \min(W^+, W^-)$.

③ Obtain the critical value, x_c , from MF19.

(this depends on: # of values that are not equal to m .
- the sample size, n and
- the level of significance, α .)

* value is different bw 1-tailed / 2-tailed.
Refer to formula booklet.

④ Reject H_0 if $x_t \leq x_c$.

$$(1) H_0: \text{median, } m = 1.8 \\ H_1: m > 1.8. \quad \begin{cases} \text{state} \\ \text{hyp} \end{cases}$$

$$(3) n = 10, \alpha = 0.05, \text{ 1-tailed}$$

$$\therefore x_c = 10$$

W_i	$ W_i - M $	sign	W^+	W^-
1.6	0.2	-		2
1.1	0.7	-		7
2.1	0.3	+	3	
2.4	0.6	+	6	
2.2	0.4	+	4	
2.9	1.1	+	10	
2.6	0.8	+	8	
2.3	0.5	+	5	
2.7	0.9	+	9	
1.9	0.1	+	1	

$$W^+ = 46 \quad W^- = 9$$

(4) Since $x_t < x_c$, we

reject H_0 .

Hence, there is sufficient evidence to suggest the population median is greater than 1.8 kg.

Normal approximation

If we take many samples from the population, and record down the values of W^+ , assuming H_0 is true,

$\Rightarrow W^+$ can be modelled as a discrete random variable.

① Since the total of the ranks = $\frac{n(n+1)}{2}$,

and (in the long run) it can be expected that $W^+ = W^- = \frac{1}{2}(\text{total})$,

$$\text{hence } E(W^+) = \frac{1}{2} \left(\frac{n(n+1)}{2} \right) = \frac{n(n+1)}{4}.$$

$$② \text{Var}(W^+) = \frac{1}{24} n(n+1)(2n+1)$$

(for now, take for granted.

It is in MFA).

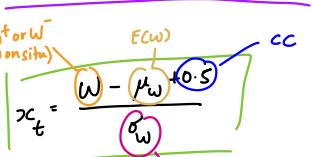
$$\Rightarrow W^+ \sim N\left(\frac{n(n+1)}{4}, \frac{n(n+1)(2n+1)}{24}\right)$$

$$\Rightarrow \text{Subsequently, } x_c = \phi^{-1}(\alpha)$$

and reject H_0 if $x_t \leq x_c$.

(see chapter 3).

and



PAIRED SAMPLE TESTS

We use paired-sample tests to compare pairs of data from 2 different populations.

↪ specifically, we compare the differences between the data.

* alternative for paired-sample t test

PAIRED SAMPLE SIGN TEST

ASSUMPTIONS

The differences bw the pairs must:

- ① be independent (from each other);
- ② come from the same population distn (median is fixed).

* the populations themselves come from "pairs" of data (eg before & after); hence they are often not independent from each other.

METHOD

The methodology is similar to the 1-sample variant; however, the sign of the difference between the values is recorded instead.

① State hypotheses $H_0: m_1 - m_2 = 0$
 $H_1: m_1 - m_2 \neq 0$

② Calculate the sign of $(P_1 - P_2)$ for each pair. ignore any Record S_+ (# of "+" signs) Values = 0. and S_- (# of "-" signs).

difference bw the two values in the pair.

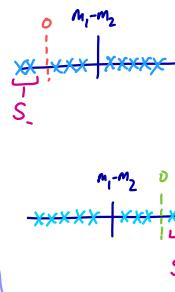
③ $S_+ \text{ or } S_- \sim B(n, \frac{1}{2})$.

Cases:

- If $H_1: m_1 - m_2 > 0$,
 $\hookrightarrow x_t = P(S_- \leq s_-)$;
- If $H_1: m_1 - m_2 < 0$,
 $\hookrightarrow x_t = P(S_+ \leq s_+)$;
- If $H_1: m_1 - m_2 \neq 0$,
 $\hookrightarrow x_t = 2P(S_+ \leq s_+)$.

④ Reject H_0 if

$$x_t < \alpha. \quad (\alpha = \text{sig level}).$$



EXAMPLE

An investigation is carried out into the effectiveness of two types of post-operative pain relief drug: Drug1 and Drug2. Seven adults agree to take Drug1 on one day, and drug 2 on the second. The time, in hours, of pain relief is recorded.

	Drug1	Drug2
A	4.1	3.9
B	3.2	3.3
C	5.3	5.0
D	5.1	4.6
E	4.2	4.6
F	3.8	3.2
G	3.6	4.3

(test, at the 5% sig level, whether drug2 gives longer pain relief than drug1).

Let $d = \text{median of drug2} - \text{median of drug1}$

(1) $H_0: d = 0$

$H_1: d > 0 \quad (m_2 > m_1)$

(2) $d_2 \quad d_1 \quad \text{sign}(d_2 - d_1)$

3.9	4.1	-
3.3	3.2	+
5.0	5.3	-
4.6	5.1	-
4.6	4.2	+
3.2	3.8	-
4.3	3.6	+

(3) $S_- \sim B(7, \frac{1}{2})$.

(a) $x_t = P(S_- \leq 4)$

$$= \left(\frac{1}{2}\right)^7 \left({}^7C_0 + {}^7C_1 + {}^7C_2 + {}^7C_3 + {}^7C_4 \right)$$

$$= 0.77343.$$

(4) $0.77343 > 0.05 \therefore x_t > \alpha$.

Hence there is sufficient evidence to accept H_0 .

Hence there is insufficient evidence to suggest that

WILCOXON MATCHED-PAIRS SIGNED-RANK TEST

ASSUMPTIONS

Again, this test has all the assumptions as the sign test, but with the additional criterion that

↳ the distribution of the differences must be symmetric.

METHOD

- ① State hypotheses. (see above)
- ② Calculate the sign of the differences (P_1, P_2), and the absolute value of the differences
- ③ Rank each difference in a similar manner to the tabulation outlined in the 1-sample test.
- ④ Obtain $x_t = \min(W_+, W_-)$.
- ⑤ From MF19, get x_c . (using $n - (\# \text{ of values} = 0)$, α and whether it is a 1 or 2-tailed test).
- ⑥ Reject H_0 if $x_t \leq x_c$.

↳ Find W_+ and W_- .

NORMAL APPROXIMATION

Like the 1-sample test, we can use a normal approximation to estimate x_t and x_c for large values of n ($\frac{n(n+1)}{2} > 20$).

$$\Rightarrow x_c = \Phi^{-1}(\alpha) \quad \text{and} \quad x_t = \frac{W - \mu_W + 0.5}{\sigma_W},$$

(or $\Phi^{-1}(1-\alpha)$)

and reject H_0 if $x_t < x_c$.

EXAMPLE (USING ABOVE Q)

$$(1) H_0 : d = 0 \quad (d = m_2 - m_1)$$

$$H_1 : d > 0.$$

$$(4) x_t = \min(12, 16) = 12.$$

$$(5) n = 7, \alpha = 0.05, \text{ 1-tailed}$$

$$\therefore x_c = 3.$$

	d_2	d_1	$d_2 - d_1$	sign	w_+	w_-	$(6) x_t > x_c$
(2)	3.9	4.1	0.2	-		2	
(3)	3.3	3.2	0.1	+	1		
	5.0	5.3	0.3	-		3	
	4.6	5.1	0.5	-		5	
	4.6	4.2	0.4	+	4		
	3.2	3.8	0.6	-		6	
	4.3	3.6	0.7	+	7		
							$\therefore w_+ = 12 \quad w_- = 16$

WILCOXON RANK-SUM TEST

The Wilcoxon rank-sum test is the non-parametric alternative to the 2-sample t-test.

↪ unlike the paired sample tests, the two samples do not have to be the same size.

METHOD

① State the hypotheses.

$$\Rightarrow H_0: m_1 - m_2 = 0 \quad (\text{let } d = m_1 - m_2).$$

$$H_1: m_1 - m_2 \begin{cases} > 0 \\ < 0 \\ \neq 0 \end{cases}$$

② Rank each sample value, taking note of the population they are from.

i.e.	P_x	P_y	w_x	w_y	rank of x & y resp.
values from pop. X	3 2 5 8	values from pop Y	3 2 4 7		
	6		5		
	1		1		
	7		6		

③ Calculate $W_x = \sum w_x$ & $W_y = \sum w_y$.

(In the above example, $W_x = 16$, $W_y = 12$).

④ Calculate the sum of the ranks of the smallest sized sample when they are ranked the "other way around", $W_{x'}$. i.e. lowest rank = highest #.

Proof of formula. Since $W_x + W_{x'} = n_x(n_x + n_y + 1)$; it follows that $W_{x'} = m(m+n+1) - W_x$.

$$W_{x'} = m(m+n+1) - W_x$$

⑤ Calculate the test statistic, x_t , where

$$x_t = \min(W_x, W_{x'}) \\ (= \min(W_y, W_{y'})) \\ (= \min(W_x, n_x(n_x + n_y + 1) - W_x)).$$

⑥ Using:

- the sample sizes of the pops n_x & n_y ,
- the SL α ,
- whether the test is 1/2-tailed;

calculate x_t .

⑦ Reject H_0 if $x_t \leq x_c$.

EXAMPLE.

Researchers are investigating the effect of vitamin B12 on the size of the brain. A sample of males aged between 25 and 40 years is selected. Nine of them are known to have low B12 levels and seven are known to have high levels. After a brain scan, the ratio of brain volume to skull capacity is recorded.

Low B12 levels	0.795	0.798	0.802	0.805	0.806	0.807	0.808	0.81	0.812
High B12 levels	0.786	0.789	0.792	0.796	0.799	0.8	0.803		

Carry out a Wilcoxon rank-sum test, at the 5% significant level, to see whether the level of vitamin B12 affects the size of the brain.

Let $d = \text{median of low B12 levels} - \text{median of high B12 levels}$.

$$(1) \Rightarrow H_0: d = 0$$

$$H_1: d \neq 0.$$

(2)

x	y	w_x	w_y
0.795	0.786	4	1
0.798	0.789	6	2
0.802	0.792	9	3
0.805	0.796	11	5
0.806	0.799	12	7
0.807	0.800	13	8
0.808	0.803	14	10
0.810		15	
0.812		16	

* you can use a pencil to take note of the ranks crossed out.

(3)

$$n_x = 9 \quad n_y = 7 \quad \therefore W_x = 100$$

$$w_y = 36$$

$$(4) W_{x'} = n_y(n_x + n_y + 1) - W_y \\ = 7(9 + 7 + 1) - 36 \rightarrow \text{in MF19.} \\ = 83.$$

$$(5) x_t = \min(W_y, W_{x'}) \\ = \min(36, 83) \\ = 36.$$

$$(6) n_x = 9 (= n), n_y = 7 (= m), \alpha = 0.05, \text{ 2-tailed} \\ \therefore x_c = 40.$$

(7) Since $x_t \leq x_c$, reject H_0 .
(make the relevant conclusions.)

NORMAL APPROXIMATION.

If n_x & n_y (n, m) are large ($n_x, n_y \geq 10$),

we can approximate x_t with a normal distribution, with $E(x_t) = \frac{m(n+m+1)}{2}$

& $\text{Var}(x_t) = \frac{mn(m+n+1)}{12}$ (comb n is the larger of n_x & n_y)

$$\Rightarrow x_t = \frac{x_t - \mu + 0.5}{\sigma} \quad \& \quad x_c = \phi^{-1}(\alpha).$$

and reject H_0 if $x_t \leq x_c$.