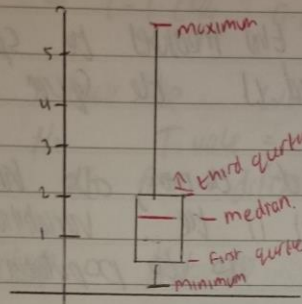I would use a boxplot to graph the 30 individual values for each method

Box plot are drawn in the following way



This is a simple way to represent statistical data on a plot in which a rectangle is draw to represent the second and third quartile (50% of data), usually with a vertical line or horizontal line inside it representing the median

The upper and lower quartiles are the vertical lines extending above and below the box

Dots or stars above or below box represent outliers

Comparing the box plots visually will give us an indication of the possible difference between the two methods

b) At $N$ = Sample Size which equals 30 for both method

Mean - For each sample the sum of the values divided by the count or value. As you can see, the means are relatively close

St.dev - Standard deviation, this measures the spread of the values from the mean and is the square root of the variance

SE Mean - Standard error of the mean, also know as Standard deviation of the mean, this is the variability of the mean of the sample compared to the population mean

Difference = mu (Type A) - mu (Type B)
This is the two sample testing the hypothesis population that there is no difference between Type A mean and type B mean population mean

If population mean are the same $U_1 - U_2 = 0$
$H_0$ : Population mean the same
$H_1$ : Population mean not the same

Estimate for difference: This is mean A - mean B which is used in the t-test formula.

95% CI for difference: Degress freedom is $n_1 + n_2 - 2 = 60 - 2 = 58$
$\alpha$ = 0.05%. Using these two value we can create a tcritical interval where t lies above and below some value.

If the calculated t lie within this interval, we fail to reject $H_0$

equation for CI i) estimate $\pm$ t critical se (error)

$$-0.569 \pm 2.008 \, (K4) \cdot 36 \qquad \frac{24}{6} \qquad \frac{25}{6}$$

$$\pm 0.7220439467$$

$$\pm 0.725$$

$$\Rightarrow (-1.294, \ 0.156)$$

T calculated value must lie within interval if we are to accept Ho. T value = -1.57 which is outside the 95% CI for difference

We can conclude that the population mean are different.

P value is also important. If too = p value is less they are of value of 0.05 we also reject null hypothesis

The p value is the probability that the mean are the same but if this probability is less than 0.05 we reject / ignore it.

∟ There are two types of errors type 1 and type 2

null hypothesis

| | true | false |
|---|---|---|
| accept | OK | Type 2 |
| reject | type 1 | OK |

Type 1 error is the significance level α which we set ourselves. Probability of rejecting the null hyp when it is true

Type 2: error is probability of accepting null hypothesis when it is false

4

## Type 1 error

The probability of a type 1 error is the level of significance as the test hypothesis denoted by alpha $\alpha$.

example: If chloesteral level of a healthy man is normally distributed with mean 180 and St. 20, and men with levels over 225 are diagnosed as not healthy, what is probability of type 1 error:

$$z = \frac{225-180}{20} = 2.25, \text{ corresponding t-tail area is } 0.0122$$

- The larger a sample size, the more likely a hypothesis test will detect a small difference
- When we are setting $\alpha = 0.05$ what we are actually saying

## Type 2 error.

Occurs when we reject alternative hypothesis when AH is true. This probability is denoted by "beta" $\beta$.

example: If by using larger values of $n$ we will give lower t critical value which will lead to a more accurate result

The "Power" test is used to calculate probability of type II.

The sample size determines the amount of sampling error inherent in a test result. Other things equal, effects are harder to detect in smaller sample sizes. Increasing sample size is easiest way to boost statistical power of the test.
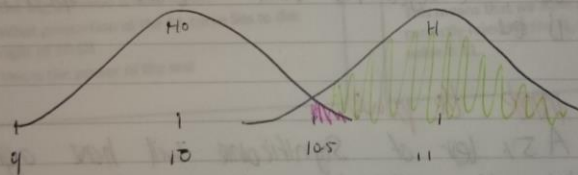
# ERRORS TYPES

We create a hypothesis
$H_0$ long term mean $= 10$
$H_1$ long term mean $\neq 10$.

We chose a different value for mean say 11.
Assume SD $= 1.5$.

What would you expect to happen when you took samples of size 21?

We look at the overlap of the sampling distributions under of $H_0$ and $H_1$ (11cm)

Create the 95% CI using t dist and N-1 df.



The two distributions overlap. But by how much?
what proportion of the 11 curve led to be right of
$(0.68/95\% \; \text{for} \; 10)$
This is the power of the test

$$ t = \frac{10.68 - 11}{0.327} = 0.9 $$

$= 83\%$

Mean we have a 83% chance of correctly rejecting
the $H_0$ when the long term value is 11

- we can increase this chance by changing
the value of n

A larger n will give a smaller S.E and thus a larger t value

larger sample size yields higher power

Statistical power is the probability of correctly rejecting a false null hypothesis when a specific alternative hypothesis is true

Power test influenced by
- Difference between actual population mean and null hypothesis mean
- variability of data $\sigma$
- Sample size N
- alpha error $\alpha$

The power against a specific alternate is calculated as the probability that the test will reject Ho when that specific alternate is true

Ways to increase the power
- Increase $\alpha$. A 5% level of significance will have a greater chance of rejecting the alternate than a 1% test because the strength of the evidence required for rejecting is less
- Consider a particular alternate that is farther away from the value of $\mu$ that are in Ho but the close to the hypothesised value $\mu$o are harder to detect that values of $\mu$ that are farm from $\mu$o
- Increase sample size. More data will provide more info about $\bar{x}$ so ve have a better chance of distinguishing values of $\mu$
- Decrease $\sigma$. This has same effect as increasing sample size, it provides more info about $\mu$ Improving the measurement process and restricting attention to a subpopulation are two common ways to decrease $\sigma$

# Stats    First    Year

Change to standardised normal    mean 0 st. of 1 by

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

For Samples   $N = \frac{\sigma}{\sqrt{n}}$ also called Standard Error
By central limit theorem

Sampling proportions

$$P = \frac{2u}{80} = 0.25$$

$$SE(p) = \sqrt{\frac{P*(1-P)}{n}} = 0.048$$

$95\%$  CI $= \hat{p} \pm 1.96 \, SE(p)$

$$0.25 \pm 1.96(0.048)$$

$n^* p^2(1-p)$  ho) to be $> 75$ to work

## Deciding how wide our CI will be

$$\hat{p} \pm 1.96^x \sqrt{\frac{P(1-P)}{n}}$$

we decide on a given error say 0.02

$$1.96^x \sqrt{\frac{P(1-P)}{n}} < 0.02$$

manipulate get us:  $n > \dfrac{\hat{p} \times (1-p) \times 1.96^2}{0.02^2}$

In worst case $p = 0.50 \Rightarrow$ sub in to get $cou$

## For Sample Sizes for estimating means

$95\%$ CI $= \bar{x} \pm 1.96^x \dfrac{Sd}{\sqrt{n}}$          d is the width

$$n \geq \frac{1.96^2 \times Sd^2}{d^2}$$

Test statistic for sample and hypothesis

$$z = \frac{\mu_0 - \text{Sample mean}}{\frac{\text{Sd}}{\sqrt{n}}}$$

$\mu_0$ is population mean

Testing population mean v sample mean

Example - population mean < 250       $\frac{250 \cdot 248.42}{\frac{8.98}{\sqrt{50}}} = 1.26$

Sample = 248.42

Sample sd Sd

8.84 = or Sample

Or we could use our z value and see if it
lie) within (-1.96, 1.96) interval, accept if within

Or calculate CI by Sample mean $\pm$ 1.96 $\times$ SE, if population mean
within CI accept $H_0$

2

<u>T-test</u>  - two sided

two group   $t = \dfrac{(\bar{x}_{new} - \bar{x}_{normal}) - hyp^{th} val}{se(\bar{x}_{new} - \bar{x}_{normal})}$

$SE[\bar{x}_1 - \bar{x}_2] = \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}$

or pool SE. $= \cdot S \times \sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}}$

$S^2 = \dfrac{(n_1 - 1) \times S_1^2 + (n_2 - 1) \times S_2^2}{n_1 + n_2 - 2}$

95% for CI : ( 2, 3)

bounce for new process is behe 2 and 3cm higher than
normal process

<u>ONE SIDED</u>

Is the population (long term mean bounce height) from new
process   higher/lower than   the usual proces

$H_0:$  $\mu_{usual} \leq \mu_{new}$       or  $\mu_{new} - \mu_{usual} \leq 0$

$H_0$  $\mu_{usual} > \mu_{new}$       or  $\mu_{new} - \mu_{usual} > 0$

$\overline{value \geq}$ critical $t$,  evidence against $H_0$

$H_0:$    $\mu$    $\geq 750$  hw

$H_1:$    $\mu$    $< 750$  hw            reject $H_0$

critical $t = $    $-1.69$

value    $< -1.69$  evidence against $H_0$

If we went $\alpha = 5\%$  we  use  two tailed but $\alpha = 10\%$
to get critical $t$

# STATS QUICK NOTES

$$p \pm 1.96 \times \sqrt{\frac{p(1-p)}{n}} \quad\quad \text{for sample proportion}$$

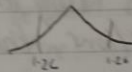$$1.96 \times \sqrt{\frac{p(1-p)}{n}} < 0.02$$

$$1.96^2 \left(\frac{p(1-p)}{n}\right) < 0.02^2$$

$$\frac{1.96^2 \, p(1-p)}{0.02^2} = n \quad\quad n = 2401$$

$$p \pm 1.96 \, \frac{st}{\sqrt{n}} \quad\quad\quad \text{for means}$$

$$z < \frac{H_0 - \text{sample pro}}{\frac{st}{\sqrt{n}}} \quad\quad t\text{-tell.} \quad \frac{H_0 - \text{sample pro}}{Se.}$$

Testing - 3 ways
1. Calculate Stand Error
2. Calculate probability of tail.



$$2 \left[ 1 - P ( z < 1.26 ) \right] = 0.2$$

3. Probability of 0.2 of getting our data
4. Compare + again α 0.05, if bigger accept $H_0$

### 2-Tail
- α = 0.05 convert to z ⇒ ±1.96
- calculate z, if in (−1.96, 1.96) accept $H_0$
  1.27 is ✓

### CI
value ± 1.96 S.e.     248.40 ± 1.96×1.27    (245.91, 250.87)
if $H_0$ value in interval accept ✓

3

# CHI-SQUARE

| Qmyra. | June | December | |
|---|---|---|---|
| Yes | 50 | 60 | 110 |
| No. | 450 | 440 | 890 |
| Total | 500 | 500 | 1000 |

Ho: no association between the two variables (whether you listen to Qmayra) and time (June vs December) in the population

i.e. Proportion of people in population who listen in June in population is the same as proportion who listen in december in population.

$$\pi_{June} = \pi_{December}$$

H₁: There is an association between the two variables, whether you listen to Qmyra and time.

i.e. Proportion of people who listen in June in population is not same as proportion who listen in December in population

$$\pi_{June} \neq \pi_{Decem}$$

Expected vote
$$\frac{row\ total \times column\ total}{N}$$

|  |  | obtne | | |
|---|---|---|---|---|
| Qmyra. | June | December | | |
| Yes | 55 | 55 | 50 | 60 |
| No. | 445 | 445 | 450 | 440 |

How To Compare
$$\frac{\sum (observed - expected)^2}{expected}$$

- will be big if observed values differ from expected

4

$x^2$ dist

$$\sum \frac{(O-E)^2}{E} = \frac{(50-55)^2}{55} + \frac{(60-55)^2}{55} + \frac{(450-445)^2}{445} + \frac{(440-445)^2}{445}$$

$$= 0.455 + 0.455 + 0.056 + 0.056 = 1.021$$

df = (number rows -1) * (number colums -1)

$(2-1) \times (2-1) = 1$

From table $\alpha = 0.05$, df = 1

critical value = 3.84

-(value) > 3.84 - evidence against $H_a$

-value $\leq$ 3.84 - not enough evidence against $H_0$

- 1.021, no evidence agant $H_0$

CI for difference in Propion

$(P_{Dec} - P_{True}) \pm t_{critial} * SE (P_{Dec} - P_{True})$

$n_1 + n_2 - 2$ df $= 50 + 50 - 2 = 98 \Rightarrow 1.96$

$$SE(P_1 - P_2) = \sqrt{\frac{P_1 \times (1-P_1)}{n_1} + \frac{P_2 \times (1-P_2)}{n_2}}$$

$$\sqrt{\frac{0.12 \times (1-0.12)}{50} + \frac{(0.10)(1-0.1)}{50}} = 0.02$$

$0.02 \pm 1.96 \cdot 0.02 = 0.02 \pm 0.04 = (-0.02; 0.06)$

0 in interval, no evidence agant difference, accept $H_0$

5.

## Chi-Square

- Chi-square does not give any info about strenght of relationship
- only conveys the existence / non existence of relationship between the variables meaningful

## Power

| | Null | hypothesis false |
|---|---|---|
| Accept | $\alpha$ True $p = 1 - \alpha$ | Type II error $p = \beta$ |
| Reject | Type 1 $p = \alpha$ | Correct $p = 1 - \beta$ |

Type II probability of accepting a false $H_0$

Power is probability of rejecting when $H_0$ is false
(1 - Type 2 error)

repeat for 11,12,13 for jack queen king

```
if (decknamed == "11") {
    decknamed [2,i] = "A"
}
```

~~compte~~ shuffled.deck.idx = smple (1:52, size=52, replar=FALSE)

shuffled.deck = decknamed [ , shuffled.deck.idx]

```
# store card.
N=1000
first.card = matrix (nrow=N, ncd=2)

for (j in ~~between~~ 1:N) {
    shuffled.deck.idx = sample (1:52, size = 52, Replac = FALSE)
    first.card = [j,] = deck.named [, shuffled.deck [.
    shuffled.deck = deck.named [, shuffled.deck.idx]
    first card [j,] = shuffled.deck
}
                    length (
no.queen = which (first.card [,] == "Queen"))
```

STATS

"\t"    tab
"\n"    new line

For (j in 1: 100)
{
    cat ("\n", j)          prints 1 to 100
}

for (j in 1:20) {
Str = paste ("Number", j)    concentenate "number" and j
print (Str)
}

x = rnorm (x, mean, sd)        rnorm (5, 10, 1) e.g

q norm - quantile -point on axis    with x% probably below it.
p norm - percentile    - cumulative
d norm - density

x = runif (100)    - uniform distribution 0,1

Gamma distribution    $\alpha$ $\beta$    $f(x) = \dfrac{\beta^{\alpha}}{\Gamma(\alpha)} x^{x-1} e^{-\beta x}$
                    shape  rate.

Set $\alpha = 1$ we get exp distribution    $f(x) = \beta e^{-\beta x}$

# shuffling a deck of cards
    deck = 1:52
Shuffled = sample (deck, size=52)    replace=FALSE    sample - shuffle anything / permute it.

Suit = c(rep('spades', 13), rep ('clubs', 13), rep ('hearts', 13), rep ('diamonds', 13))
    # gives each suit named out 13 times

denominator = c (1:13, 1:13, 1:13, 1:13)
deck named = rbind (Suit, denomination)    rbind for columns

For (decknamed i in 1:52)

if ( decknamed [2, i] == "1") {
        decknamed [2, i] = "ace"    second row > the number
}