

Stat 135 Lec 21 (no lec 20)

Remaining 2 quizzes pushed back 1 week

Quiz 3: Nov 8

Quiz 4: Nov 22

Last time

sec 11.2 Comparing the means of independent normal populations with the same variance.

$$X_1, \dots, X_n \sim N(\mu_x, \sigma^2), Y_1, \dots, Y_m \sim N(\mu_y, \sigma^2)$$

σ^2 known:

$$\bar{X} - \bar{Y} \sim N(\mu_x - \mu_y, \sigma^2(\frac{1}{n} + \frac{1}{m}))$$

2 sample Z-test.

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0, 1)$$

σ^2 unknown:

we use unbiased estimator of σ^2 :

Pooled sample variance

$$S_p^2 = \frac{(n-1)S_x^2 + (m-1)S_y^2}{n+m-2}$$

(2 sample t-test) $t = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$

if $n+m-2 \geq 30$ you can use 2 sample Z-test.

- Today
- Sec 11.2
- { (1) Duality of HT and CI.
 - (2) Comparing means of independent normal populations with different variances)
 - (3) Power calculations >
- Sec 11.3
- (3) Comparing Pooled Samples

Post Midterm topics

- parametric test (Normal theory) (sec 11.2, 11.3)
t-test) comparing the means of
2 normal populations
- non parametric test (sec 11.2, 11.3)
comparing the means of
2 continuous populations.
- Chi-square tests testing whether 2 or more categorical variables have the same distribution.
- ANOVA (Normal theory) (sec 12.2)
F-test) comparing the means of
2 or more normal populations
- Regression (chap 14)
- Bayesian Statistics (parameter estimation)
hypothesis testing

ex (Problem 11.6.10) (Duality of HT and CI)

Verifying that the two sample t-test at level α of $H_0: \mu_x - \mu_y = 0$

$$H_1: \mu_x - \mu_y \neq 0$$

accepts iff the $100(1-\alpha)\%$ CI for $\mu_x - \mu_y$ contains zero.

Soln

accept null if

$$\left| \frac{\bar{x} - \bar{y}}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \right| < t_{\frac{n+m-2}{2}, \alpha/2}$$

$$\Leftrightarrow -t(\frac{\alpha}{2}) < \frac{\bar{x} - \bar{y}}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} < t(\frac{\alpha}{2})$$

$$\Leftrightarrow -t(\frac{\alpha}{2}) S_p \sqrt{\frac{1}{n} + \frac{1}{m}} < \bar{x} - \bar{y} < t(\frac{\alpha}{2}) S_p \sqrt{\frac{1}{n} + \frac{1}{m}}$$

$$\Leftrightarrow \bar{x} - \bar{y} - t(\frac{\alpha}{2}) S_p \sqrt{\frac{1}{n} + \frac{1}{m}} < 0 < \bar{x} - \bar{y} + t(\frac{\alpha}{2}) S_p \sqrt{\frac{1}{n} + \frac{1}{m}}$$

$$\Leftrightarrow 0 \in \bar{x} - \bar{y} \pm t(\frac{\alpha}{2}) S_p \sqrt{\frac{1}{n} + \frac{1}{m}}$$

Sec 11.2

We have assumed X, Y both have same variance

If we remove this assumption

then $\frac{S_x^2}{n} + \frac{S_y^2}{m}$ is a natural

unbiased estimator of $\text{Var}(\bar{x} - \bar{y})$

$$\text{chck } E\left(\frac{S_x^2}{n} + \frac{S_y^2}{m}\right) = \frac{E(S_x^2)}{n} + \frac{E(S_y^2)}{m}$$

$$= \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m} = \text{Var}(\bar{x}) + \text{Var}(\bar{y})$$

$$= \text{Var}(\bar{x} - \bar{y})$$

$\swarrow X, Y \text{ indep}$

$\Rightarrow \frac{S_x^2}{n} + \frac{S_y^2}{m}$ is unbiased estimator of $\text{Var}(\bar{x} - \bar{y})$.

but what is

$$\frac{\bar{x} - \bar{y} - 0}{\sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}} ?$$

Ans it is approximately a t distribution w/ df given by formula on p428.

P428 Rice

We have used the assumption that the two populations have the same variance. If the two variances are not assumed to be equal, a natural estimate of $\text{Var}(\bar{X} - \bar{Y})$ is

$$\frac{s_X^2}{n} + \frac{s_Y^2}{m}$$

If this estimate is used in the denominator of the t statistic, the distribution of that statistic is no longer the t distribution. But it has been shown that its distribution can be closely approximated by the t distribution with degrees of freedom calculated in the following way and then rounded to the nearest integer:

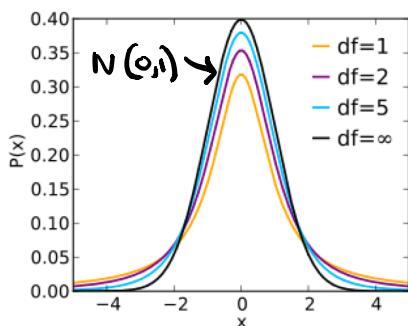
$$df = \frac{\left[\left(s_X^2/n\right) + \left(s_Y^2/m\right)\right]^2}{\frac{\left(s_X^2/n\right)^2}{n-1} + \frac{\left(s_Y^2/m\right)^2}{m-1}}$$

↗ no need
to memorize,

need to calculate
this in R
for hw,

≤ 50 rats were randomly divided into 2 groups of 25 each. The rats in one grp were given steroids. They were timed running a maze. The rats in the non steroid grp had an avg time of 10 seconds, with SD of 2 seconds. The rats in the steroid grp had an avg time of 9 seconds and SD=5.

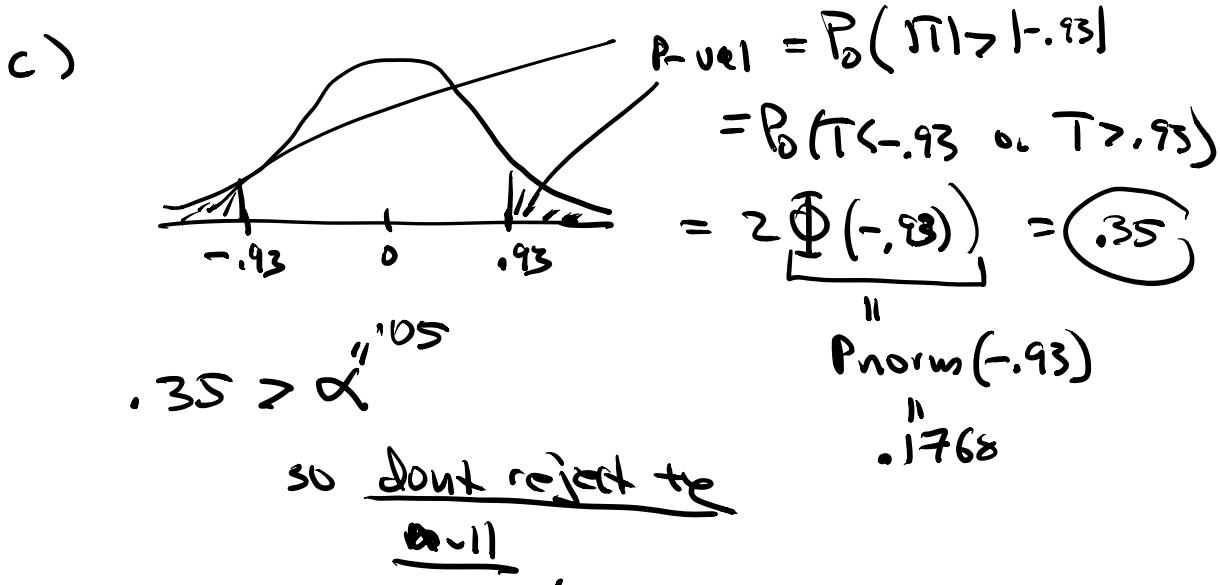
- What is the SE of differences of the 2 averages?
- Calculate a test statistic for testing the null hypothesis that steroids make no difference.
- What is the p-value of the test?
(using formula for df of t-test gives 31 df. This is large so just do a z-test).



$$a) S_{\bar{x}_T - \bar{x}_C} = \sqrt{\frac{s_T^2}{n_T} + \frac{s_C^2}{n_C}} = 1.08$$

$$b) H_0: \mu_T = \mu_C \quad T = \frac{\bar{x}_T - \bar{x}_C - 0}{1.08} = -0.93$$

$$H_1: \mu_T \neq \mu_C$$



Sec 11.2.2 Power

Calculations of power are an important part of planning an experiment in order to determine how large of a sample size you should use.

Let assume σ known

σ same in both populations

$$n = m$$

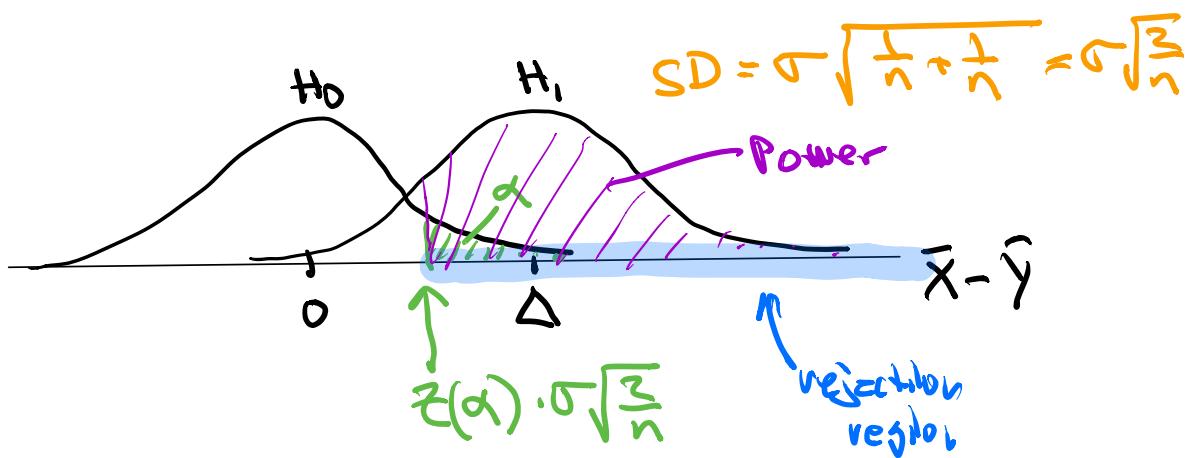
Picture

$$H_0: \mu_x - \mu_y = 0$$

$$H_1: \mu_x - \mu_y = \Delta > 0$$

one sided alternative.

α level of significance.



$$\begin{aligned}
 \text{Power} &= P_{H_1}(\alpha(x)=1) \\
 &= P_{H_1}(\bar{x}-\bar{y} > z(\alpha) \sigma \sqrt{\frac{2}{n}}) \\
 &= P_{H_1}\left(\frac{\bar{x}-\bar{y}-\Delta}{\sigma \sqrt{\frac{2}{n}}} > \frac{z(\alpha) \sigma \sqrt{\frac{2}{n}} - \Delta}{\sigma \sqrt{\frac{2}{n}}}\right) \\
 &= \boxed{1 - \Phi\left(z(\alpha) - \frac{\Delta}{\sigma \sqrt{\frac{2}{n}}}\right)}
 \end{aligned}$$

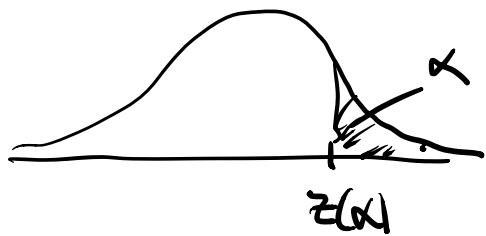
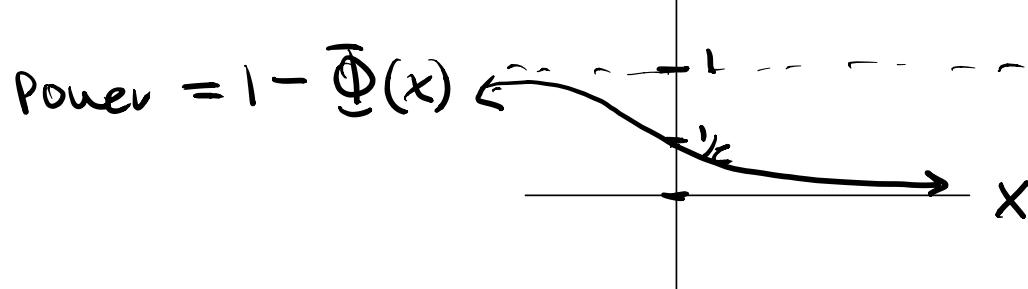
As increase $\Delta, \alpha, \sigma, n$
what happens to power?

$\Delta \uparrow$ Power \uparrow

$\alpha \uparrow$ Power \uparrow

$\sigma \uparrow$ Power \downarrow

$n \uparrow$ Power \uparrow



Ex Suppose n measurements are taken from a treatment group and independently n are taken from a control group. An SD of a single measurement is $\sigma = 10$ for both groups.

How large should n be so that the test

$$H_0: \mu_x - \mu_y = 0$$

$$H_1: \mu_x - \mu_y = \Delta > 0$$

has power .7 if $\Delta = \mu_x - \mu_y = 1$ and $\alpha = .05$?

$$\text{Power} = 1 - \Phi\left(z(\alpha) - \frac{\Delta}{\sigma}\sqrt{\frac{n}{2}}\right)$$

$$\Rightarrow \Phi\left(z(\alpha) - \frac{\Delta}{\sigma}\sqrt{\frac{n}{2}}\right) = 1 - \text{Power} = .3$$

$$q_{\text{norm}}(.3) = -.524$$

$$\Rightarrow z(\alpha) - \frac{\Delta}{\sigma}\sqrt{\frac{n}{2}} = -.524$$

$$\alpha = .05 \Rightarrow z(\alpha) = 1.645$$

$$\Delta = 1$$

$$\sigma = 10$$

$$\Rightarrow \sqrt{\frac{n}{2}} = (.524 + 1.645) \frac{10}{1}$$

$$\Rightarrow n = 941$$

Stat 135lec 22

Last time

2-sample t-test:

Compared means of independent normal populations
with different variances

$$S_{\bar{x}-\bar{y}} = \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{n}} \quad \text{standard error.}$$

Today

Sec 11.3 Paired 2 sample t-test

Paired and unpaired tests in R

Pre 9.5 The multinomial distribution and
Chi-squared test

Sec 11.3.1 paired t-test.

We wish to estimate $\mu_x - \mu_y$ but in paired studies our samples are no longer independent.

$$\text{let } x_1, \dots, x_n \text{ iid } N(\mu_x, \sigma_x^2)$$

$$y_1, \dots, y_n \text{ iid } N(\mu_y, \sigma_y^2)$$

$$\sigma_{xy} = \text{Cov}(x_i, y_j) \neq 0.$$

$$\text{let } D_i = x_i - y_i$$

$$E(D_i) = \mu_x - \mu_y \quad \text{call this } \mu_D$$

$$\text{Var}(D_i) = \text{Var}(x_i - y_i) = \sigma_x^2 + \sigma_y^2 - 2\sigma_{xy}$$

$$\bar{D} \sim N(\mu_x - \mu_y, \frac{1}{n}(\sigma_x^2 + \sigma_y^2 - 2\sigma_{xy}))$$

Notice If x, y are pos correlated

$$\sigma_{xy} > 0 \quad \text{and} \quad \text{Var}(\bar{D}) < \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{n} \quad \leftarrow \text{Var}(\bar{x} - \bar{y})$$

$$\Rightarrow \text{efficiency} = \frac{S_{\text{paired}}}{S_{\text{unpaired}}} < 1 \quad \begin{matrix} \text{when } x, y \\ \text{indep.} \end{matrix}$$

\Rightarrow larger t.s. \Rightarrow so more likely to reject the null for paired data.

$$t = \frac{\bar{D} - \mu_D}{S_{\bar{D}}} \sim t_{n-1}$$

$$H_0: \mu_D = 0 \quad \text{has rejection region } |\bar{D}| > t_{n-1}(\alpha/2) S_{\bar{D}}$$

$$H_1: \mu_D \neq 0$$

lec 22 in-class exercise on b-courses

Stat135 lecture 22

Start Over

Examine your data

```
X <- c(24.6,17,16,10.4,8.2,7.9,8.2,7.9,5.8,5.4,5.1,4.7)
Y <- c(10.1,5.7,5.6,3.4,6.5,0.7,6.5,0.7,6.1,4.7,2.0,2.9)
n <- length(X)
df <- data.frame(X,Y)
df
```

| | X <dbl> | Y <dbl> |
|--|------------|------------|
| | 24.6 | 10.1 |
| | 17.0 | 5.7 |
| | 16.0 | 5.6 |
| | 10.4 | 3.4 |
| | 8.2 | 6.5 |
| | 7.9 | 0.7 |
| | 8.2 | 6.5 |
| | 7.9 | 0.7 |
| | 5.8 | 6.1 |
| | 5.4 | 4.7 |

1-10 of 12 rows

Previous 1 2 Next

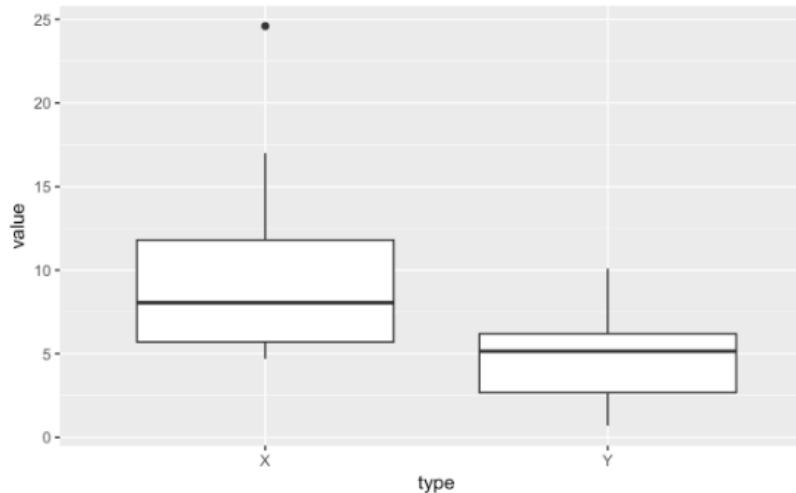
reorganize data so that it is "tidy"

```
library(tidyr)
#we reformat the data frame to tidy format for easier analysis
df_tidy<- df %>% gather(key=type, value=value, `X`, `Y`)
head(df_tidy)
```

| | type <chr> | value <dbl> |
|---|---------------|----------------|
| 1 | X | 24.6 |
| 2 | X | 17.0 |
| 3 | X | 16.0 |
| 4 | X | 10.4 |
| 5 | X | 8.2 |
| 6 | X | 7.9 |

6 rows

```
library(ggplot2)
df_tidy %>% ggplot(aes(x=type,y=value)) + geom_boxplot()
```



```
df_tidy
```

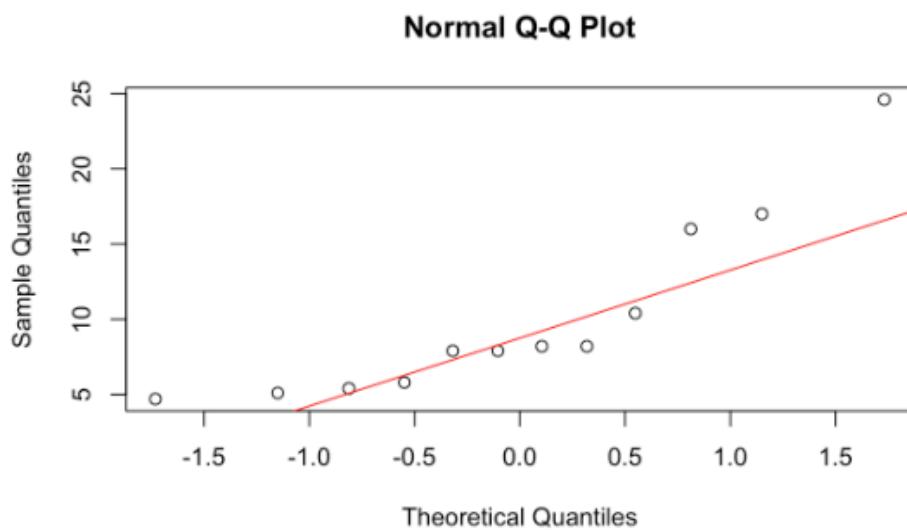
| type | value |
|------|-------|
| X | 24.6 |
| X | 17.0 |
| X | 16.0 |
| X | 10.4 |
| X | 8.2 |
| X | 7.9 |
| X | 8.2 |
| X | 7.9 |
| X | 5.8 |
| X | 5.4 |

1-10 of 24 rows

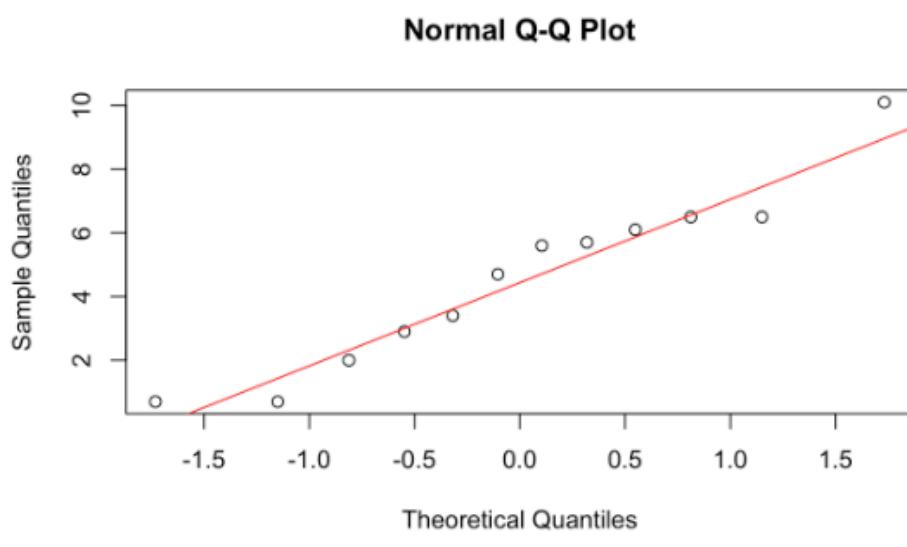
Previous 1 2 3 Next

Lets examine that the data is approximately normal and that the variances are approximately equal.

```
qqnorm(X); qqline(X,col=2)
```



```
qqnorm(Y); qqline(Y,col=2)
```



```
S.X <- sd(X) #sd() in R is sample SD  
S.Y <- sd(Y)  
S.X
```

```
## [1] 5.80517
```

```
S.Y
```

```
## [1] 2.636957
```

Neither appear to be very normal and the variances don't appear to be the same. At this point I would do a nonparametric test called the Mann-Whitney test (after midterm). But lets do a parametric test and assume X and Y are normal but not assume the variances are the same.

Test that two independent samples come from populations with the same mean.

Since the variances aren't the same lets use the formula for degrees of freedom in the t-test on page 428 rounded to the nearest integer. Here n=m.

$$\text{deg.freedom} = \frac{((S_X^2/n) + (S_Y^2/n))^2}{(S_X^2/n)^2/(n-1) + (S_Y^2/n)^2/(n-1)}$$

```
deg.freedom <- round((n-1)*(S.X^2/n+S.Y^2/n)^2/((S.X^2/n)^2+(S.Y^2/n)^2))  
deg.freedom
```

```
## [1] 15
```

$H_0: \mu_X - \mu_Y = 0$

$H_1: \mu_X - \mu_Y \neq 0$

$\alpha = .01$

```
t <- (mean(X)-mean(Y))/sqrt(S.X^2/n+S.Y^2/n)  
t
```

```
## [1] 3.001744
```

```
p_val <- 2*(1-pt(abs(t),df=deg.freedom))  
p_val
```

```
## [1] 0.008940922
```

p.val < .01 so **reject** null.

Test that two paired samples come from populations with the same mean.

We assume that before and after are normal vectors.

```
before <- c(24.6,17,16,10.4,8.2,7.9,8.2,7.9,5.8,5.4,5.1,4.7)
after <- c(10.1,5.7,5.6,3.4,6.5,0.7,6.5,0.7,6.1,4.7,2.0,2.9)
D <- before-after
S.D <- sd(D)
```

```
D<- before-after
t <- (mean(D))/(S.D/sqrt(n))
t
```

```
## [1] 4.179141
```

bigger t value since SE for paired
→ smaller than for unpaired

```
p_val <- 2*(1-pt(abs(t),df=n-1))
p_val
```

```
## [1] 0.001538782
```

smaller than p-val for unpaired test
since t value bigger,
This is good if alternative true, bad otherwise.

$p_{\text{val}} < .01$ so reject the null. Notice that SE for the sampling distribution for the paired test is less than the sampling distribution for the unpaired test. In other words paired sampling is more efficient.

```
(SE_pair <- S.D/sqrt(n) ) #SE of paired sampling distribution
```

```
## [1] 1.322042
```

```
(SE_unpair <- sqrt((S.X^2/n)^2+(S.Y^2/n)^2)) #SE of unpaired sampling distribution
```

```
## [1] 2.867492
```

```
SE_pair/SE_unpair #efficiency
```

```
## [1] 0.4610447 < 1,
```

Sec 9.5 The multinomial distribution

Generalizes the binomial distribution

The multinomial RV is a sum of n independent outcome (cells) RV.

$m=2$ for binomial,

$$\left. \begin{array}{l} P_1 = \text{Prob of outcome 1} \\ P_2 = " " " 2 \\ \vdots \\ P_m = \text{Prob of outcome } m \end{array} \right\} P_1 + \dots + P_m = 1$$

What is chance of getting x_1 outcome 1, x_2 outcome 2, ... x_m outcome m ? ($n = x_1 + \dots + x_m$)

Ans $\binom{n}{x_1, x_2, \dots, x_m} P_1^{x_1} \dots P_m^{x_m}$ multinomial formula,
 $\frac{n!}{x_1! \dots x_m!}$

If we know P_1, \dots, P_m the multinomial distribution
is completely specified.

(i.e. we can calculate the prob of any outcome)

In n trials

$$\text{Obs} \quad \overbrace{\boxed{x_1 | x_2 | \dots | x_m}}^{\sum_{i=1}^m x_i! = n}$$

$$\text{Exp} \quad \boxed{np_1 | np_2 | \dots | np_m} \quad \sum_{j=1}^m p_j = 1$$

$$x_i \sim \text{Bin}(n, p_i) \approx N(np_i, np_i q_i)$$

Chi square Statistic

$$\sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} \xrightarrow{n \rightarrow \infty} \chi^2_{m-1}$$

When we have 2 outcomes (binomial)
we will show this is χ_1^2

$$\sum_{i=1}^2 \frac{(O_i - E_i)^2}{E_i} = \frac{(X_1 - np_1)^2}{np_1} + \frac{(X_2 - np_2)^2}{np_2}$$

argue $\frac{(X_1 - np_1)^2}{np_1 q_1} = z^2 \sim \chi_1^2$

$np_1 q_1$
 $\text{Var}(k_1)$

more generally

$$\sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} = \sum_{i=1}^{m-1} \frac{(O_i - E_i)^2}{E_i q_i} \xrightarrow{\text{for } \lambda = 4e^{-n}} \sum_{i=1}^{m-1} z_i^2 \sim \chi_{m-1}^2$$

algebra

QK Please review the example of modeling the arrival of particles in a Poisson random scatter in lecture 3 (intro to chp 8).

We show the table below which is the observed and expected number of arrivals in 1207 10 second intervals.

We model the number of arrivals in a 10 second interval as $X \sim \text{Pois}(\lambda \cdot 10)$. We estimated $\hat{\lambda} = .84$ so $\lambda \cdot 10 = 8.4$ arrivals/10 sec.

$$\text{Here } P_1 = e^{-8.4} + \frac{e^{-8.4} (8.4)}{1!} + \frac{e^{-8.4} (8.4)^2}{2!}$$

$$P_2 = \frac{e^{-8.4} (8.4)^3}{3!} = .022$$

etc.

Expected count for $n=3$ variables is $1207(.022)=27$

Think of a box with 16 outcomes with probabilities P_1, P_2, \dots, P_{16}

$$\underbrace{[P_1, P_2, P_3, \dots, P_{16}]}_{\downarrow 1207 draws}$$

| | n | Observed | Expected |
|----------|-----|----------|----------|
| P_1 | 0-2 | 18 | 12.2 |
| P_2 | 3 | 28 | 27.0 |
| P_3 | 4 | 56 | 56.5 |
| . | 5 | 105 | 94.9 |
| . | 6 | 126 | 132.7 |
| . | 7 | 146 | 159.1 |
| . | 8 | 164 | 166.9 |
| . | 9 | 161 | 155.6 |
| . | 10 | 123 | 130.6 |
| . | 11 | 101 | 99.7 |
| . | 12 | 74 | 69.7 |
| . | 13 | 53 | 45.0 |
| . | 14 | 23 | 27.0 |
| . | 15 | 15 | 15.1 |
| . | 16 | 9 | 7.9 |
| P_{16} | 17+ | 5 | 7.1 |

1207 1207

Null: Our observed counts come from the multinomial distribution $MN(1207, P_1(\lambda), P_2(\lambda), \dots, P_{16}(\lambda))$ where $P_i(\lambda)$ is Poisson(λ).

Alt: Our observed counts come from some other multinomial distribution

O | | | | | | | |

E | T T T T C C

these boxes
are the table

we calculate χ^2 stat

$$\sum_{\text{all rows}} \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{\underbrace{16-1}_{\substack{\uparrow \\ \text{dim of} \\ \text{Sample} \\ \text{space}}} - 1}$$

dim of null space

(P_1, \dots, P_{16})

Where

$\sum_{i=1}^{16} P_i = 1$

$(\lambda \text{ is 1 unknown parameter in } \text{Pois}(\lambda))$

Last time

Sec 9.5 goodness of fit χ^2 test.

In chap 11 we did 1,2 sample + test for the mean of a normal RV. Our box had a continuous distribution w/ mean M , var σ^2 .

Now we have a categorical RV having m outcomes having probabilities P_1, \dots, P_m . The chance a sample of size n for our box has a certain composition is given by the multinomial

$$\text{formula } P(X_1=x_1, X_2=x_2, \dots, X_m=x_m) = \frac{n!}{x_1! x_2! \dots x_m!} P_1^{x_1} \dots P_m^{x_m}$$

The goal is to test whether a model for the population distribution (P_1, \dots, P_m) fits our data.

We draw n times with replacement from our box and get observed counts $[x_1 | x_2 | \dots | x_m], x_1 + \dots + x_m = n$

If the prob of tickets in box is $P_1(\theta), \dots, P_m(\theta)$ we expect to get counts $[nP_1(\theta) | nP_2(\theta) | \dots | nP_m(\theta)]$

We do a goodness of fit test.

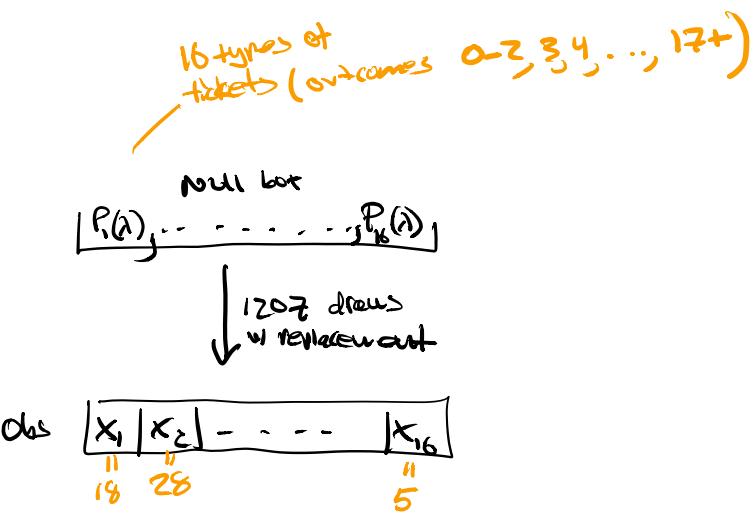
Pearson Chi-Square T.S. $\sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{m-1-k}$ where k is the dimension of θ .

A goodness of fit test explores how good your probability model fits your data,

If the p-value of our Pearson Chi-Square T.S. is smaller than α we reject the null hypothesis.

ex

| n | Observed | Expected |
|-----|----------|----------|
| 0-2 | 18 | 12.2 |
| 3 | 28 | 27.0 |
| 4 | 56 | 56.5 |
| 5 | 105 | 94.9 |
| 6 | 126 | 132.7 |
| 7 | 146 | 159.1 |
| 8 | 164 | 166.9 |
| 9 | 161 | 155.6 |
| 10 | 123 | 130.6 |
| 11 | 101 | 99.7 |
| 12 | 74 | 69.7 |
| 13 | 53 | 45.0 |
| 14 | 23 | 27.0 |
| 15 | 15 | 15.1 |
| 16 | 9 | 7.9 |
| 17+ | 5 | 7.1 |
| | 1207 | 1207 |



Null: Our observed counts come from the multinomial distribution $MN(1207, P_1(\lambda), P_2(\lambda), \dots, P_{16}(\lambda))$ where $P_i(\lambda) \sim \text{Poisson}(\lambda)$.

Alt: Our observed counts come from some other multinomial distribution

we calculate χ^2 stat

$$\sum_{\text{all rows}} \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{M-1 - K}$$

$\underbrace{16-1}_{M-1} - 1$

Compare P value to $\alpha = .05$.

Today

Sec 9.5 GLRT for the multinomial distribution,

① Examples

② theory: We will see $-2 \log \Lambda \approx \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i}$

we know $-2 \log \Lambda \sim \chi^2_{\dim \mathcal{L} - \dim \mathcal{W}_0}$

$$\Rightarrow \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{\dim \mathcal{L} - \dim \mathcal{W}_0}$$

↑ same space, null space, same spec.

① Hardy Weinberg (HW) equilibrium model

If gene frequencies are in equilibrium then
 genotypes AA, Aa, aa occur in the population
 w prob $(1-\theta)^2$, $2\theta(1-\theta)$, θ^2 according to the
 HW model.

$$\text{Null box} \quad \begin{cases} P_1(\theta) = (1-\theta)^2 \\ P_2(\theta) = 2\theta(1-\theta) \\ P_3(\theta) = \theta^2 \end{cases}$$

We observe blood type

| | AA | Aa | aa |
|----------|----------------|----------------|----------------|
| M | 342 | 500 | 187 |
| Obj | x ₁ | x ₂ | x ₃ |
| n = 1029 | | | |

Step 1 Find $\hat{\theta}_{ML}$

$$l(\theta) = \frac{n!}{x_1! x_2! x_3!} (1-\theta)^{2x_1} (2\theta(1-\theta))^{x_2} \theta^{2x_3}$$

$$l(\theta) = \log n! - \sum_{i=1}^3 \log x_i! + x_1 \log (1-\theta)^2 + x_2 \log 2\theta(1-\theta) + x_3 \log \theta^2$$

$$l'(\theta) = 0 \Rightarrow \hat{\theta}_{ML} = \frac{2x_3 + x_2}{2n} = \frac{2(187) + 500}{2(1029)} = \boxed{.4277}$$

$$\therefore P_1(\hat{\theta}) = (1 - .4277)^2 = \boxed{.3200}$$

$$P_2(\hat{\theta}) = 2(.4277)(1 - .4277) = \boxed{.4887}$$

$$P_3(\hat{\theta}) = (.4277)^2 = \boxed{.1804}$$

Step 2

| | | | |
|-----|----------------------------|--------|--------|
| Exp | $n\hat{P}_i(\hat{\theta})$ | 502.83 | 185.60 |
| " | | 340.57 | |

Test

H_0 : have $MN(1029, .32, .49, .18)$

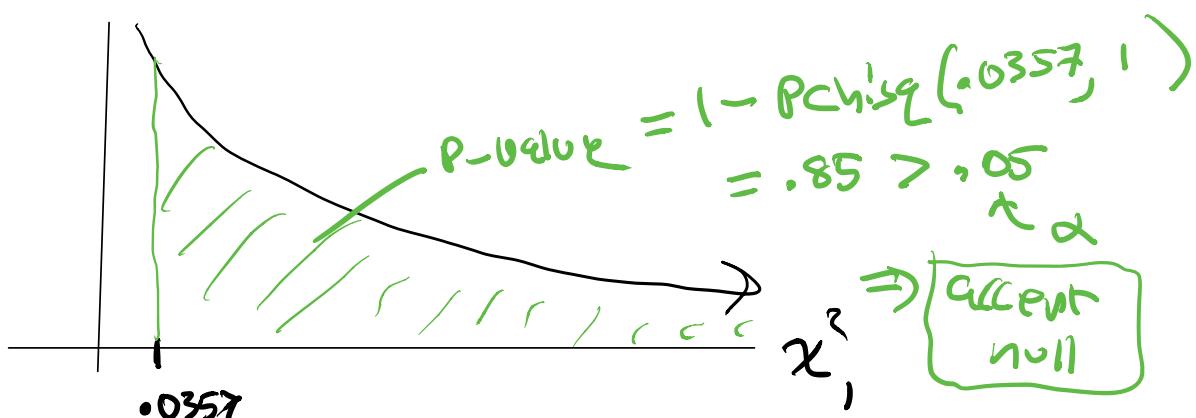
H_1 : have some other $MN \sim n=1029$.

| | | | |
|---|-----|-----|-----|
| O | 342 | 500 | 187 |
|---|-----|-----|-----|

| | | | |
|---|-------|-------|-------|
| E | 340.6 | 502.8 | 183.6 |
|---|-------|-------|-------|

$$\chi^2 = \frac{(342 - 340.6)^2}{340.6} + \dots$$

$$= .0357$$



② Sec 9.5 Theory.

H_0 : have a MN distribution of cell prob

$$P(\Theta) = (P_1(\Theta), \dots, P_m(\Theta)) \quad \text{where } \Theta = (\Theta_1, \dots, \Theta_k) \in \Omega_0 \subseteq \mathbb{R}^k$$

open interval

H_1 : have a MN dist w/ different

cell prob than our null,

Ω = set of m non-negative numbers that sum to 1
 $\underbrace{(P_1, \dots, P_m)}_{\text{w/ } P_1 + P_2 + \dots + P_m = 1}$

Sample Space.

$$\Lambda = \frac{\max_{\Theta \in \Omega_0} (\text{lik}(P_1(\Theta), \dots, P_m(\Theta)))}{\max_{(P_1, \dots, P_m) \in \Omega} (\text{lik}(P_1, \dots, P_m))}$$

\Downarrow

$$\text{lik}(\hat{P}_1, \dots, \hat{P}_m)$$

$$\text{Fact } \hat{P}_i = \frac{x_i}{n} \quad \leftarrow \text{see Sec 8.5.1}$$

P273
use Lagrange Multipliers.

$$x_i = n\hat{P}_i$$

By consistency of MLE

$$P_i(\hat{\theta}) \rightarrow P_i(\theta)$$

and $\hat{P}_i \rightarrow P_i(\theta)$ if null is true

$$\begin{aligned} \Lambda &= \frac{\frac{n!}{x_1! \cdots x_m!} P_1(\hat{\theta})^{x_1} \cdots P_m(\hat{\theta})^{x_m}}{\frac{n!}{x_1! \cdots x_m!} \hat{P}_1^{x_1} \cdots \hat{P}_m^{x_m}} \\ &= \prod_{i=1}^m \left(\frac{P_i(\hat{\theta})}{\hat{P}_i} \right)^{x_i} \\ -2 \log \Lambda &= -2 \sum_{i=1}^m (n \hat{P}_i) \log \left(\frac{n P_i(\hat{\theta})}{n \hat{P}_i} \right) \\ &= \boxed{2 \sum_{i=1}^m O_i \log \left(\frac{O_i}{E_i} \right)} \end{aligned}$$

Fact (See P342 (Taylor series argument))

$$2 \sum_{i=1}^m O_i \log \left(\frac{O_i}{E_i} \right) \approx \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i}$$

K Pearson χ^2 statistic

$$\dim \mathcal{R} = m-1$$

$\omega_0 = \left\{ (\theta_1, \dots, \theta_k) \text{ s.t. } \theta_i \text{ is in an open interval in } \mathbb{R} \right\}$

$$\dim \omega_0 = k.$$

$$\Rightarrow \begin{cases} -2 \log 1 \rightarrow \chi^2_{m-1-k} \\ \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} \end{cases}$$

Last time

Sec 9.5 Application of GLRT to goodness of fit χ^2 test.

The goodness of fit (G.O.F) χ^2 test tests whether a categorical random variable follows a certain distribution (null hypothesis), versus that it follows some other distributions.

We see that $-2 \log \lambda$, where λ is from the GLRT (see theory in lec 23), is approximately the Pearson χ^2 test statistic, # cells

$$\sum_{i=1}^{\text{# cells}} \frac{(O_i - E_i)^2}{E_i}.$$

Since $-2 \log \lambda \sim \chi^2_{(\text{dim sample space} - \text{dim null space})}$, it follows that

$$\sum_{i=1}^{\text{# cells}} \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{(\text{dim sample space} - \text{dim null space})}.$$

Today

(1) Assumptions of goodness of fit test,

(2) Examples.

①

Assumptions of goodness of fit test

- ② have one categorical variable
- ③ have independent observations
(drawn w/ replacement)
- ④ The outcomes are mutually exclusive,
(i.e. every ticket goes in exactly one cell)
- ⑤ we require large n and no more than 20% of expected counts < 5

Small E_i and large n then

$$E_i = n p_i \Rightarrow p_i \text{ small}$$

$$X_i \sim \text{Bin}(n, p_i)$$

want

$$X_i \sim N(n p_i, n p_i z_i)$$

but for small p_i , $X_i \sim \text{Poisson} (\mu = n p_i)$
not normal,

② Examples of GOF test

ex As part of a study on the selection of grand juries in Alameda county, the educational level of grand jurors was compared with the county distribution:¹⁴

| <i>Educational level</i> | <i>County</i> | <i>Number of jurors</i> |
|--------------------------|---------------|-------------------------|
| Elementary | 28.4% | 1 |
| Secondary | 48.5% | 10 |
| Some college | 11.9% | 16 |
| College degree | 11.2% | 35 |
| | | — |
| Total | 100.0% | 62 |

Could a simple random sample of 62 people from the county show a distribution of educational level so different from the county-wide one?

Solution

- As part of a study on the selection of grand juries in Alameda county, the educational level of grand jurors was compared with the county distribution:¹⁴

| <i>Educational level</i> | <i>County</i> | <i>Number of jurors</i> |
|--------------------------|---------------------------------|-------------------------|
| Elementary | $P_1 = 28.4\% \quad E_1 = nP_1$ | 1 |
| Secondary | $P_2 = 48.5\% \quad E_2 = nP_2$ | 10 |
| Some college | $P_3 = 11.9\% \quad E_3 = nP_3$ | 16 |
| College degree | $P_4 = 11.2\% \quad E_4 = nP_4$ | 35 |
| Total | 100.0% | 62 = n |

Could a simple random sample of 62 people from the county show a distribution of educational level so different from the county-wide one?

$$H_0 : \text{multinom}(62, P_1 = 28.4, P_2 = 48.5, P_3 = 11.9, P_4 = 11.2)$$

$$H_1 : \text{some other multinomial dist.}$$

| O | G |
|----|------|
| 1 | 17.6 |
| 10 | 29.8 |
| 16 | 7.4 |
| 35 | 6.9 |

$$\chi^2_3 = \frac{(1-17.6)^2}{17.6} + \frac{(10-29.8)^2}{29.8}$$

$$+ \frac{(16-7.4)^2}{7.4} + \frac{(35-6.9)^2}{6.9}$$

$$= 153.24$$

$$\chi^2_3 (.05) = 7.814$$

reject null

reject null

Ex

(10 pts) Suppose the school claims that the number of incoming students for the past 4 years has been steadily increasing by 10% per year. Assuming that there is no dropout, form a null and alternative hypothesis around this claim. Name the test to be carried out, calculate the degrees of freedom, and list out the observed values (O) and their corresponding expected values (E).

Note: You do NOT need to calculate the final result or perform the test.

try and solve this.

-) (10 pts) Suppose the school claims that the number of incoming students for the past 4 years has been steadily increasing by 10% per year. Assuming that there is no dropout, form a null and alternative hypothesis around this claim. Name the test to be carried out, calculate the degrees of freedom, and list out the observed values (O) and their corresponding expected values (E).

Note: You do NOT need to calculate the final result or perform the test.

Solution:

We perform a Goodness of Fit test.

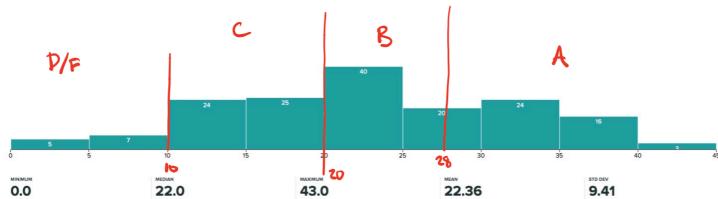
H_0 : The number of students in each year is increasing by 10% per year.

H_1 : The number of students in each year is not as described.

Since we are assuming that there are no students dropping out, the null hypothesis states that the number of students in each year in our sample follows a multinomial distribution with probabilities equal to the proportion of students in each year. To calculate the exact null distribution, suppose the proportion of seniors in Berkeley is p , then under the null the probabilities of a randomly selected student being a Freshman, Sophomore, Junior and Senior are $(p * 1.1^3, p * 1.1^2, p * 1.1, p)$, respectively. Under the constraint that the sum of these 4 probabilities is 1, we can solve for p : $p \approx 0.215$. Plugging in p and the total sample size ($n = 904$), the expected numbers of students in each year are $904 * (p * 1.1^3, p * 1.1^2, p * 1.1, p) = (259.26, 235.69, 214.26, 194.79)$ and our observed values are $(241, 227, 216, 220)$. The degrees of freedom for this test is $4 - 1 = 3$

Next time, test of homogeneity and independence,
Sec 13.3, 13.4

Stat 155 Lec 25



Nice job on a
tough midterm!

Last time sec 9.5 Goodness of fit,

Test that one (categorical) variable has certain multinomial distribution with cell prob $P_1(\theta), \dots, P_m(\theta)$ where $\theta = (\theta_1, \dots, \theta_k)$ and $P_1(\theta) + \dots + P_m(\theta) = 1$.

Key result: $\sum_{\text{# cells}} \frac{(O_i - E_i)^2}{E_i} \approx \chi^2_{\dim S - \dim w_0}$

\uparrow \uparrow
 $m-1$ K

Today

Sec 13.3, 13.4 Test of Homogeneity and Independence

1 cat variable

GOF does the cat var. have the claimed multivariate distribution?

test of homogeneity

- do 2 or more subgrps of a pop share the same multivariate dist?

2 cat variables — test of indep — are the 2 cat variables independent?

Sec 13.3 Test of homogeneity

Ex 13.2.1

A study was done comparing frequencies of a particular allele in a sample of diabetics and non diabetics.

Data was observed:

| | | Diab | non Diab | tot |
|---|----------|------|----------|-----|
| I | Bb or bb | 12 | 4 | |
| | BB | 39 | 49 | 88 |
| | | 51 | 53 | 104 |

Are the relative frequency of the alleles significantly different in the 2 groups Diab, non Diab

Theory

J indep, I cell multinomials

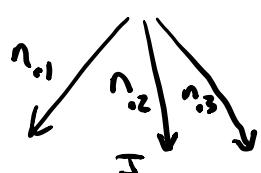
H_0 : all J multinomials are same

H_1 : not all to same

Pictorial

$$\begin{aligned} I &= 2 \\ J &= 3 \end{aligned}$$

$$\text{null bal} \quad \frac{\pi_1}{\pi_2} = \frac{\hat{\pi}_1}{\hat{\pi}_2} = \frac{\pi_1}{\pi_2}$$



$$\begin{aligned} \hat{\pi}_1 &= \frac{n_{11}}{n} \\ \hat{\pi}_2 &= \frac{n_{21}}{n} \end{aligned}$$

obs

I

| | | | |
|----------|----------|----------|----------|
| n_{11} | n_{12} | n_{13} | n_{10} |
| n_{21} | n_{22} | n_{23} | n_{20} |

$$\begin{array}{cccc} n_{11} & n_{12} & n_{13} & n_{10} \\ n_{21} & n_{22} & n_{23} & n_{20} \\ \hline n_{11} & n_{12} & n_{13} & n \\ \hline \end{array}$$

Exp

| | | |
|----------------------------|----------------------------|----------------------------|
| $\frac{n_{1,1}n_{1,1}}{n}$ | $\frac{n_{1,2}n_{1,2}}{n}$ | $\frac{n_{1,3}n_{1,3}}{n}$ |
| $\frac{n_{2,1}n_{2,1}}{n}$ | $\frac{n_{2,2}n_{2,2}}{n}$ | $\frac{n_{2,3}n_{2,3}}{n}$ |

What if we divide population according to allele type instead of according to whether or not you have diabetes? You get the same Obs and Exp tables!!.

Notice if you have I index, J cell multinomials

Null Box

| | |
|-----------------------|-----------------------------|
| π_1, π_2, π_3 | $n_{1,1}, n_{1,2}, n_{1,3}$ |
| π_1, π_2, π_3 | $n_{2,1}, n_{2,2}, n_{2,3}$ |

Obs

| | | | |
|-----------|-----------|-----------|-----------|
| $n_{1,1}$ | $n_{1,2}$ | $n_{1,3}$ | $n_{1,1}$ |
| $n_{2,1}$ | $n_{2,2}$ | $n_{2,3}$ | $n_{2,1}$ |
| $n_{1,1}$ | $n_{1,2}$ | $n_{1,3}$ | n |

$\hat{\pi}_1 = \frac{n_{1,1}}{n}$

$\hat{\pi}_2 = \frac{n_{1,2}}{n}$

$\hat{\pi}_3 = \frac{n_{1,3}}{n}$

$n_{1,1}$ ← Same as before

$n_{2,1}$ ← Same as before

Exp

| | | |
|--------------------------------|--------------------------------|--------------------------------|
| $\frac{n_{1,1}\hat{\pi}_1}{n}$ | $\frac{n_{1,2}\hat{\pi}_2}{n}$ | $\frac{n_{1,3}\hat{\pi}_3}{n}$ |
| $\frac{n_{2,1}\hat{\pi}_1}{n}$ | $\frac{n_{2,2}\hat{\pi}_2}{n}$ | $\frac{n_{2,3}\hat{\pi}_3}{n}$ |

$\hat{\pi}_1 = \frac{n_{1,1}}{n}$ ← Same as before

$$\chi^2_{df} = \dim \mathcal{L} - \dim w_0$$

w_0 = have 1 multinomial (null box)
w/ $I-1$ free param

$$\dim w_0 = I-1$$

\mathcal{L} = have J multinomials in group (a)
w/ $I-1$ free params per multinomial
 $\Rightarrow \dim \mathcal{L} = J(I-1)$

$$\Rightarrow df = \dim \mathcal{L} - \dim w_0 = J(I-1) - (I-1)$$

$\boxed{-(I-1)(J-1)}$

T.S

$$\chi^2_{(I-1)(J-1)} = \sum_{i=1}^I \sum_{j=1}^J \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$= \sum_{i,j} \left(\frac{(n_{ij} - \frac{n_i n_j}{n})^2}{\frac{n_i n_j}{n}} \right)$$

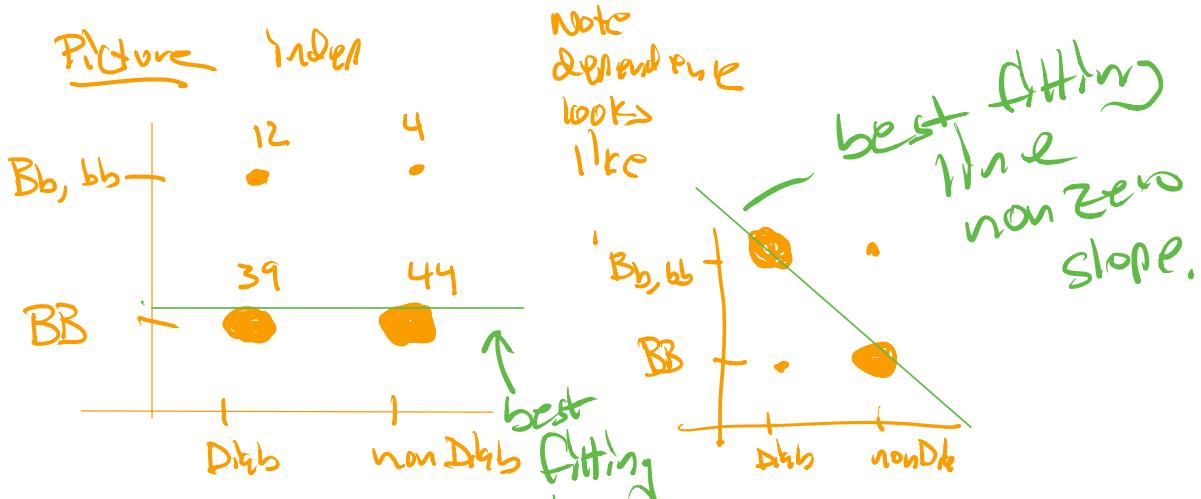
Sec 13.4 Test of Independence

- test whether 2 categorical variables are independent.

Ex modify Problem 13.8.1

Same data but diff question.

Is there a relationship between whether you have diphtheria and your allele type?



Theory Slope 0

1 J cell multinomial

1 I cell multinomial

Notation

$$P(J=j) = \overline{\pi}_{J=j} = \pi_j$$

$$P(I=i) = \overline{\pi}_{I=i} = \pi_i$$

$$P(I=i, J=j) = \overline{\pi}_{I=i, J=j} = \pi_{ij}$$

Independence means $\pi_{ij} = \pi_i \pi_j$

H_0 : two RV are indep (i.e. $\pi_{ij} \subset \pi_i \pi_j$)

H_{alt}

H_1 : two RV are dep (i.e. $\pi_{ij} \neq \pi_i \pi_j$ some i, j)

J

Picture

I=2

J=3 I

| | π_1 | π_2 | π_3 |
|---------|---------------|---------------|---------------|
| π_1 | $\pi_1 \pi_1$ | $\pi_1 \pi_2$ | $\pi_1 \pi_3$ |
| π_2 | $\pi_2 \pi_1$ | $\pi_2 \pi_2$ | $\pi_2 \pi_3$ |

null box



Obs

$$\begin{bmatrix} n_{11} & n_{12} & n_{13} \\ n_{21} & n_{22} & n_{23} \\ n_{01} & n_{02} & n_{03} \end{bmatrix} \begin{array}{l} n_{1..} \\ n_{2..} \\ n_{0..} \end{array}$$

$$\hat{\pi}_{I=1} = \frac{n_{1..}}{n}$$

$$\hat{\pi}_{J=1} = \frac{n_{0..}}{n}$$

Exp

$$\left[\begin{matrix} \frac{n_{1,1}}{n} & \frac{n_{1,2}}{n} & \frac{n_{1,3}}{n} \\ \frac{n_{2,1}}{n} & \frac{n_{2,2}}{n} & \frac{n_{2,3}}{n} \end{matrix} \right]$$

Same as before.

Same T.S ||

We will show next time that
 $\text{df} = (I-1)(J-1)$ so test of
 independence and test of
 homogeneity are equivalent tests
 for 2 different experimental
 designs.

Start 13.5 Dec 26

Last time Sec 9.5 Sec 13.3, 13.4

- - -

Goodness of Fit — single sample from population

Answers Q: Does population have a particular multinomial distribution?

Test of Homogeneity — multiple samples from subgrps
of a single population

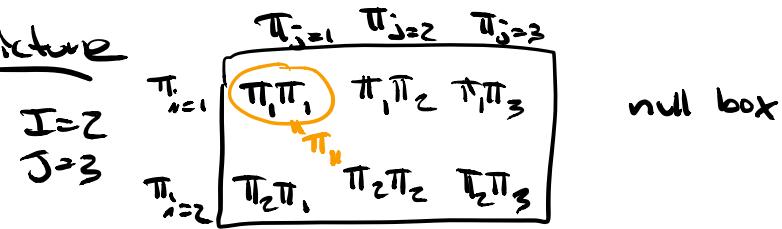
Answers Q: Do subgrps have same multinomial dist?

Test of Independence — single sample from population

Answers Q: Are two categorical variables independent?

Test of independence H_0 : two cat variables indep
 H_1 : dependent.

Picture



$\downarrow n$

Obs

$$\begin{bmatrix} n_{11} & n_{12} & n_{13} \\ n_{21} & n_{22} & n_{23} \\ n_{11} & n_{12} & n_{13} \end{bmatrix} \quad \begin{matrix} n_{10} \\ n_{20} \\ n_{00} \end{matrix}$$

Some tables are
in test of
Homogeneity.

Exp

$$\begin{bmatrix} \frac{n_{11}\pi_{11}}{n} & \frac{n_{12}\pi_{12}}{n} & \frac{n_{13}\pi_{13}}{n} \\ \frac{n_{21}\pi_{21}}{n} & \frac{n_{22}\pi_{22}}{n} & \frac{n_{23}\pi_{23}}{n} \end{bmatrix}$$

Today

- ① d.f. of Test of independence (Sec 13.4)
- ② in class exercise
- ③ non parametric test version of unpaired t-test (Sec 11.2.3)
"The Mann Whitney Test".

① d.f. & test of Independence (Sec 13.4)

$w_0 = \dim$ of marginal prob of null box

$$\pi_{ij} \quad i=1,..,I \quad \pi_{ij} \quad j=1,..,J$$

or constraints $\sum_{i=1}^I \pi_{ij} = 1, \quad \sum_{j=1}^J \pi_{ij} = 1$

$$\Rightarrow \dim w_0 = I-1 + J-1$$

$$\mathcal{L} = \left\{ \pi_{ij} \quad i=1,..,I \quad j=1,..,J \quad \text{s.t.} \quad \sum_{i=1}^I \pi_{ij} = 1 \right.$$

$$\dim \mathcal{L} = IJ-1$$

$$\Rightarrow \text{d.f.} = \dim \mathcal{L} - \dim w_0$$

$$= IJ-1 - ((I-1) + (J-1)) = \boxed{(I-1)(J-1)}$$

Summary

Test of Homo and Indep have the same T.S.,
and d.f. and differ by design of experiment.

In test of Indep, observational units are collected at random from the box and 2 categorical variables are observed for each unit.

In test of Homo, the data is collected by random sampling from each subgroup of the subpopulation separately,

Assumptions of χ^2 test

- ① have one or two categorical variables
- ② have Indep observations (!, e.g. draw w/o replacement)
Note SRS or if relatively small sample size.
- ③ Outcomes are mutually exclusive (every ticket is exactly one cell)
- ④ We require n to be large and no more than 20% of expected counts < 5 .

n large since $-2\log \Lambda \rightarrow \chi^2_{\dim \mathcal{L} - \dim w_0}$

if E_i is small and n is large

$$E_i = n p_i \Rightarrow p_i \text{ small}$$

$$X_i \sim \text{Bin}(n, p_i)$$

$$\text{Want } X_i \sim N(np_i, np_i \bar{p}_i)$$

but if p_i is too small, $X_i \sim \text{Poisson}(\mu = np_i)$
not normal.

② In-class exercise

For each question answer

- 1) What test to use
- 2) Null and alternative hypothesis
- 3) Are assumptions met
- 4) Degrees of freedom
- 5) Solve using chisq.test() in R

A)

9. This problem considers some more data on Jane Austen and her imitator (Morton 1978). The following table gives the relative frequency of the word *a* preceded by (PB) and not preceded by (NPB) the word *such*, the word *and* followed by (FB) or not followed by (NFB) *I*, and the word *the* preceded by and not preceded by *on*.

| Words | <i>Sense and Sensibility</i> | <i>Emma</i> | <i>Sanditon I</i> | <i>Sanditon II</i> |
|-------------------|------------------------------|-------------|-------------------|--------------------|
| <i>a PB such</i> | 14 | 16 | 8 | 2 |
| <i>a NPB such</i> | 133 | 180 | 93 | 81 |
| <i>and FB I</i> | 12 | 14 | 12 | 1 |
| <i>and NFB I</i> | 241 | 285 | 139 | 153 |
| <i>the PB on</i> | 11 | 6 | 8 | 17 |
| <i>the NPB on</i> | 259 | 265 | 221 | 204 |

Was Austen consistent in these habits of style from one work to another? Did her imitator successfully copy this aspect of her style?

B)

16. A market research team conducted a survey to investigate the relationship of personality to attitude toward small cars. A sample of 250 adults in a metropolitan area were asked to fill out a 16-item self-perception questionnaire, on the basis of which they were classified into three types: cautious conservative, middle-of-the-roader, and confident explorer. They were then asked to give their overall opinion of small cars: favorable, neutral, or unfavorable. Is there a relationship between personality type and attitude toward small cars? If so, what is the nature of the relationship?

| Attitude | Personality Type | | |
|-------------|------------------|---------|----------|
| | Cautious | Midroad | Explorer |
| Favorable | 79 | 58 | 49 |
| Neutral | 10 | 8 | 9 |
| Unfavorable | 10 | 34 | 42 |

C)

To test whether a die is fair, someone rolls it 600 times. On each roll, he just records whether the result was even or odd, and large (4, 5, 6) or small (1, 2, 3). The observed frequencies turn out as follows:

| | <i>Large</i> | <i>Small</i> |
|------|--------------|--------------|
| Even | 183 | 113 |
| Odd | 88 | 216 |

Question: Is the die fair?

A

Test of Homogeneity since independent sample from each subgroup of population. There are a couple parts to the question the first one is whether Austin was consistent in her style. For this we will analyze the first 3 columns of the table to see whether they have the same multinomial distribution.

Null: each of the 3 multinomial distributions are the same

Alt: Not all the multinomial distributions are the same

Assumptions are met:

- 1) Each ticket is multinomial with one cell checked off
- 2) Not clear if this is a random sample of frequency of words but any case ok.
- 3) check expectation table

First we set up the data:

```
data <- c(14 , 16 , 8 , 133 , 180 , 93, 12 , 14 , 12 , 241 , 285 , 139 , 11 , 6 , 8, 259 , 265 , 221)
data.matrix <- matrix(data, nrow = 6, ncol = 3, byrow = TRUE)
data.matrix

##      [,1] [,2] [,3]
## [1,]    14    16     8
## [2,]   133   180    93
## [3,]    12    14    12
## [4,]   241   285   139
## [5,]    11     6     8
## [6,]   259   265   221
```

Now perform the test (first displaying the expected table for fun):

```
chisq.test(data.matrix, correct = FALSE)$expected

##      [,1]      [,2]      [,3]
## [1,] 13.281168 15.184142  9.534690
## [2,] 141.898800 162.230569 101.870631
## [3,] 13.281168 15.184142  9.534690
## [4,] 232.420449 265.722483 166.857068
## [5,]  8.737611  9.989567  6.272822
## [6,] 260.380803 297.689098 186.930099

chisq.test(data.matrix, correct = FALSE)

## 
## Pearson's Chi-squared test
##
## data: data.matrix
## X-squared = 23.287, df = 10, p-value = 0.009735
```

Thus, we reject the null; Austen was not consistent among her works.

To see if Austen's imitator was successful in imitating her style, we can first compile all of Austen's data into one vector (otherwise it isn't quite clear which work to compare the imitator to since Austin wasn't consistent in her habits of style). We compare with Sanditon 2 as an example:

```
austen <- c()          # initialize an empty vector
for(i in 1:6){
  austen[i] = sum(data.matrix[i, ])
}
imitator <- c(2, 81, 1, 153, 17, 204)
matrix <- matrix(c(austen, imitator), nrow = 6, ncol = 2, byrow = FALSE)
chisq.test(matrix, correct = FALSE)

## 
## Pearson's Chi-squared test
##
## data: matrix
## X-squared = 29.726, df = 5, p-value = 1.67e-05
```

Thus, once again, we reject the null; the imitator was unsuccessful in imitating Austen's style.

B)

Test of independence since a single SRS of 250.

Null: two categorical variables personality type and attitude are independent

Alt: dependent

Assumptions:

- 1) every ticket is the cross classification of two multinomials with one cell checked off.
- 2) Every ticket is independent
- 3) check expectation

```
b <- matrix(c(79,58,49,10,8,9,10,34,42),nrow=3,ncol=3,byrow = TRUE)
b
```

```
##      [,1] [,2] [,3]
## [1,]    79    58    49
## [2,]    10     8     9
## [3,]    10    34    42
```

```
chisq.test(b)$expect
```

```
##      [,1]     [,2]     [,3]
## [1,] 61.585284 62.20736 62.20736
## [2,] 8.939799 9.03010 9.03010
## [3,] 28.474916 28.76254 28.76254
```

df=2*2=4

```
chisq.test(b)
```

```
##
##  Pearson's Chi-squared test
##
## data: b
## X-squared = 27.289, df = 4, p-value = 1.737e-05
```

we reject the null that variables are independent.

C)

Tricky, looks like test of independence but it is a goodness of fit test since test of independence doesn't answer the question about fairness. We have

Evan and Large= 4,6 — 183 Evan and Small= 2 — 113 Odd and Large=5 — 88 Odd and Small= 1,3 — 216

Expect probabilities 2/6,1/6,1/6,2/6

Null: fair die

Alt: unfair die

Assumptions:

- 1) Every ticket is a multinomial with 4 cells with one cell checked off
- 2) tickets are independent since rolling die is indep
- 3) Like drawing with replacement from a box
- 4) expected table:

```
c <- c(183,113,88,216)
c
```

```
## [1] 183 113 88 216
```

```
chisq.test(c,p=c(2/6,1/6,1/6,2/6))$expect
```

```
## [1] 200 100 100 200
```

df=3

```
chisq.test(c,p=c(2/6,1/6,1/6,2/6))
```

```
##
##  Chi-squared test for given probabilities
##
## data: c
## X-squared = 5.855, df = 3, p-value = 0.1189
```

is accept
die is not fair, we ~~reject~~ the null.

(3) Mann-Whitney Test (Sec 11.2.3)

In section 11.2 (unpaired 2 sample t-test)

we are comparing 2 independent samples
of normal data,

If normality doesn't hold for small sample size

use nonparametric test, called

Mann-Whitney test
Rank sum test
Wilcoxon test in R } All the same test.

lec 27

QB Friday 60F sec 9.5, Homogeneity, Independence sec 13.3
13.4

Today

sec 11.2 Mann-Whitney test — non parametric test
a.k.a. Wilcoxon - Rank Sum test version of unpaired t-test,

Motivation

We won't assume normality of our data.

Notation

X small grp $\sim F$
(usually control)
 x_1, \dots, x_n

Y larger grp $\sim G$
(usually treatment grp)
 y_1, \dots, y_m (independent samples)

ex

$$\begin{array}{ll} x_1 = 1000 & y_1 = 1400 \\ x_2 = 1380 & y_2 = 1600 \\ x_3 = 1200 & y_3 = 1180 \\ & y_4 = 1220 \end{array}$$

Sort $\{x_1, \dots, x_n, y_1, \dots, y_m\}$ from smallest to largest

$$1000, 1180, 1200, 1220, 1380, 1400, 1600$$

1 2 3 4 5 6 7

T_x, T_y are rank sums of x, y .

$$\text{ex } T_x = 9, T_y = 19$$

Null $F = G$ (compare with $\mu_x = \mu_y$ in t-test)

Alt $F \neq G$

If Null true, T_x or T_y shouldn't be too big or small

we test whether a diuretic works, for $\alpha=0.1$, one sided test.
 Two groups $X = \text{control}$ $Y = \text{treatment}$
 2 independent samples X_1, X_2, X_3 Y_1, Y_2, Y_3, Y_4

Null \rightarrow treatment doesn't work
 So the control and treatment populations are the same.

Alt treatment work (so you pee more).

We rank daily urine production from smallest = 1 to largest = 7

We find rank sum of X (smaller grp)

$$T_X = 9 \quad \leftarrow \text{rank sum of grp } X$$

How extreme is this?

Find sampling distribution of T_X .

There are $\binom{7}{3} = 35$ equally likely rank sums of X .

We see that $T_X = 9$ has

one side p-value of $\frac{7}{35} = .2$

The cutoff for $\alpha=0.1$ is $T_X = 7$ and $9 > 7$ so accept the null,

see distribution of T_X on next page

see table on next page

Example of Wilcoxon Rank Sum test from Glantz's Primer of Biostatistics

| Placebo (Control) = X | | Drug (Treatment) = Y | |
|-------------------------------|-------|-------------------------------|-------|
| Daily Urine Production (mL/d) | Rank* | Daily Urine Production (mL/d) | Rank* |
| 1000 | 1 | 1400 | 6 |
| 1380 | 5 | 1600 | 7 |
| 1200 | 3 | 1180 | 2 |
| | | 1220 | 4 |
| $T = 9$ | | | |

*1 = smallest; 7 = largest.

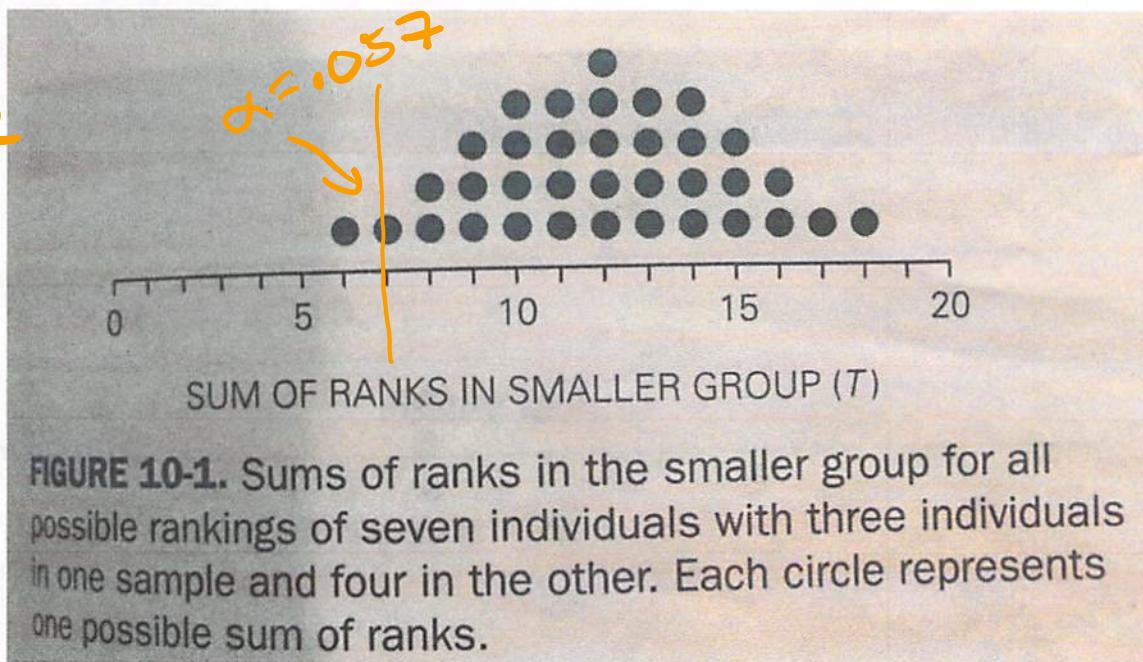


TABLE 8 Critical Values of Smaller Rank Sum for the Wilcoxon Mann-Whitney Test

| n_2 | α for Two-Sided Test | α for One-Sided Test | n_1 (Smaller Sample) | | | | | | | | | |
|-------|-----------------------------|-----------------------------|------------------------|-----|----|----|----|----|----|----|----|----|
| | | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| 3 | .20 | .10 | | 3 | 7 | | | | | | | |
| | .10 | .05 | | | 6 | | | | | | | |
| | .05 | .025 | | | | | | | | | | |
| | .01 | .005 | | | | | | | | | | |
| 4 | .20 | .10 | | 3 | 7 | 13 | | | | | | |
| | .10 | .05 | | | 6 | 11 | | | | | | |
| | .05 | .025 | | | | 10 | | | | | | |
| | .01 | .005 | | | | | | | | | | |
| 5 | .20 | .10 | 4 | 8 | 14 | 20 | | | | | | |
| | .10 | .05 | 3 | 7 | 12 | 19 | | | | | | |
| | .05 | .025 | | 6 | 11 | 17 | | | | | | |
| | .01 | .005 | | | | 15 | | | | | | |
| 6 | .20 | .10 | 4 | 9 | 15 | 22 | 30 | | | | | |
| | .10 | .05 | 3 | 8 | 13 | 20 | 28 | | | | | |
| | .05 | .025 | | 7 | 12 | 18 | 26 | | | | | |
| | .01 | .005 | | | 10 | 16 | 23 | | | | | |
| 7 | .20 | .10 | 4 | 10 | 16 | 23 | 32 | 41 | | | | |
| | .10 | .05 | 3 | 8 | 14 | 21 | 29 | 39 | | | | |
| | .05 | .025 | | 7 | 13 | 20 | 27 | 36 | | | | |
| | .01 | .005 | | | 10 | 16 | 24 | 32 | | | | |
| 8 | .20 | .10 | 5 | 11 | 17 | 25 | 34 | 44 | 55 | | | |
| | .10 | .05 | 4 | 9 | 15 | 23 | 31 | 41 | 51 | | | |
| | .05 | .025 | 3 | 8 | 14 | 21 | 29 | 38 | 49 | | | |
| | .01 | .005 | | | 11 | 17 | 25 | 34 | 43 | | | |
| 9 | .20 | .10 | 1 | 5 | 11 | 19 | 27 | 36 | 46 | 58 | 70 | |
| | .10 | .05 | 4 | *10 | 16 | 24 | 33 | 43 | 54 | 66 | | |
| | .05 | .025 | 3 | 8 | 14 | 22 | 31 | 40 | 51 | 62 | | |
| | .01 | .005 | | 6 | 11 | 18 | 26 | 35 | 45 | 56 | | |

(continued)

We could have tried to show that T_y is extremely large instead of T_x being extremely small. T_x and T_y are related.

$$T_x + T_y = 1 + 2 + 3 + \dots + 7 = \frac{7 \cdot 8}{2}$$

We used T_x in this problem since it is the smaller rank sum and the table above gives the critical value for the smaller rank sum.

Convention!
 Two populations $X \sim F$ $Y \sim G$
 Two independent samples x_1, \dots, x_n y_1, \dots, y_m

You look at T_y not T_x .

Wilcoxon tried to estimate $\pi = P(X < Y)$ and in so doing rediscovered the Mann-Whitney test.

It is useful to show how he did it because many people don't use T_x , or T_y to test H_0 but rather the Wilcoxon T.S.

To approximate $\pi = P(X < Y)$, Wilcoxon

$$\text{uses } \hat{\pi} = \frac{1}{mn} \sum_{i=1}^n \sum_{j=1}^m I(x_i < y_j)$$

$\sum_{j=1}^m I(x_i < y_j)$ called Wilcoxon T.S.
 (called W in R).

The $I(X_i \leq Y_j)$ for large $n+m$ are almost independent so CLT works well and U_Y is almost normal.

We will show that $X_1 - Y_2, X_2 - Y_2$ are dependent RVs, but $X_1 - Y_2, X_2 - Y_2$ are indep.

$$U_Y = T_Y - \frac{m(m+1)}{2}$$

so,

$$T_Y = U_Y + \frac{m(m+1)}{2} \stackrel{\text{approx}}{\sim} N\left(\frac{m(n+n+1)}{2}, \frac{mn(m+n+1)}{12}\right)$$

by theorem A above

How does U_Y relate to T_Y ?

let $X_{(1)}, \dots, X_{(n)}$ be the order

statistics of X_1, \dots, X_n

$$\hat{x}_1 = 2, \hat{x}_2 = -1, \hat{x}_3 \approx 0$$

$$X_{(1)} = -1, X_{(2)} = 0, X_{(3)} = 2$$

Note $\sum_i \sum_j I(X_i \leq Y_j)$ is +ve

Sum of indicators of all comparisons
of X and Y so

$$\sum_i \sum_j I(X_{(i)} < Y_{(j)}) = \sum_i \sum_j I(X_{(i)} < Y_{(j)})$$

Let $R_{Y_1} =$ the rank of $Y_{(1)}$ in the combined sample
 $\therefore R_{Y_1} = \#(X' < Y_{(1)}) + 1$

$R_{Y_2} =$ the rank of $Y_{(2)}$ in the combined sample
 $\therefore R_{Y_2} = \#(X' < Y_{(2)}) + 2$

etc

We can write

$$U_Y = \sum_{i=1}^n \sum_{j=1}^m I(X_{(i)} < Y_{(j)})$$

$$= (\text{number of } X' < Y_{(1)})^{R_{Y_1}-1}$$

$$+ (\text{number of } X' < Y_{(2)})^{R_{Y_2}-2} + \dots$$

$$+ (\text{number of } X' < Y_{(m)})^{R_{Y_m}-m}$$

$$U_Y = \sum_{i=1}^m R_{Y(i)} - (1+2+\dots+m)$$

"

$\frac{m(m+1)}{2}$

Sum of rank of $Y_{(i)}$ = Sum of rank of Y_i
"
 T_Y

$\Rightarrow U_Y = T_Y - \frac{m(m+1)}{2}$

Fact

Thm A P438

If $F = G$ (null)

$$E(T_Y) = \frac{m(m+n+1)}{2}$$

$$\text{Var}(T_Y) = \frac{mn(m+n+1)}{12}$$

Pt/ easy see book.

U_Y normal $\Rightarrow T_Y$ normal 

also $T_Y \stackrel{\text{approx}}{\sim} N\left(\frac{m(m+n+1)}{2}, \frac{mn(m+n+1)}{12}\right)$

$$U_Y \sim N\left(\frac{mn}{2}, \frac{mn(m+n+1)}{12}\right)$$

Called Wilcoxon W
T.S. in R.

Last time

Sec 11.2, 3 The Mann Whitney Test. (nonparametric test)

Notation

X small grp
(usually control)
 x_1, \dots, x_n

Y larger grp
(usually treatment grp)
 y_1, \dots, y_m
(Independent samples)

T_x, T_y are rank sums of X, Y .

Null $F = G$

Alt $F \neq G$

If null true expect T_x not to be extreme.
Table 8 p A21 gives critical values of T_x .

$$\text{Here } T_x + T_y = 1+2+3+\dots+(m+n) = \frac{(m+n)(m+n+1)}{2}$$

For large $m+n$, assuming the Null is true,

$$T_y \stackrel{\text{approx}}{\sim} N\left(\frac{m(m+n+1)}{2}, \frac{mn(m+n+1)}{12}\right)$$

The Wilcoxon T.S. is

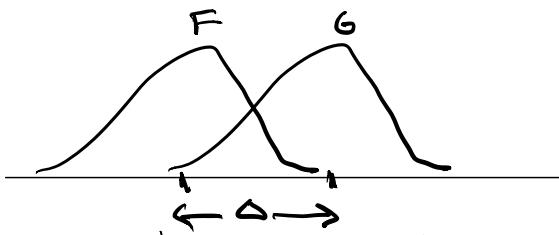
$$U_y = \sum_{i=1}^m \sum_{j=1}^n I(x_i < y_j) \stackrel{\text{assuming the null is true}}{\sim} N\left(\frac{mn}{2}, \frac{mn(m+n+1)}{12}\right)$$

these are related by $U_y = T_y - \frac{m(m+1)}{2}$ Wilcoxon T.S. (called W) in R.

Today

① Sec 11.2.3 Wilcoxon Rank Sum test in R

Given 2 independent samples x_1, \dots, x_n and y_1, \dots, y_m
`wilcox.test()` finds the T.S. W (this is U_y)
If willing to assume F and G are shifts of
one another (i.e. $G(x) = F(x - \Delta)$)



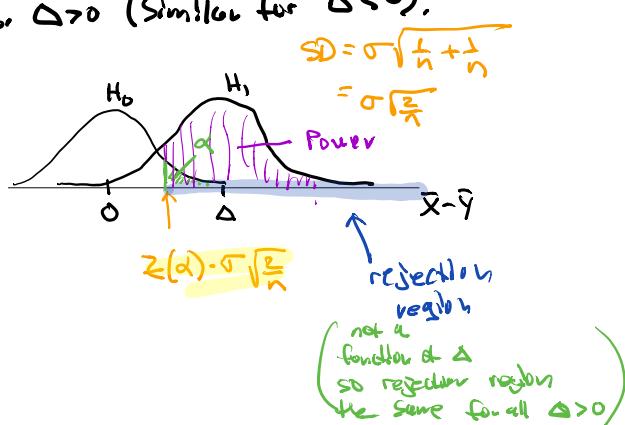
then `wilcox.test` provides a 95% CI of
what Δ is.

Note t-test is uniformly most powerful for $\Delta \neq 0$

$$H_0: M_x - M_y = 0$$

$$H_1: M_x - M_y = \Delta \neq 0$$

Picture for $\Delta > 0$ (similar for $\Delta < 0$).



So if assumptions for t-test are satisfied,
it is more powerful than non-parametric test

② Sec 11.3.2 The Signed Rank Test

① Sec 11.2.3 Wilcoxon Rank Sum test in R In-class exercise.

Unpaired 2 sample tests

✓ a: (data cleaning)

This problem is #21b in Rice chapter 11 exercises.

A study was done to compare the performance of engine bearings made of different compounds. Ten bearings of each type were tested. The table below gives the times until failure.

Discuss with your neighbor the code below:

```
library(tidyverse)
#time to failure for two types of engine bearings
type1 <- c(3.03, 5.53, 5.60, 9.30, 9.92, 12.51, 12.95, 15.21, 16.04, 16.84)
type2 <- c(3.19, 4.26, 4.47, 4.53, 4.67, 4.69, 12.78, 6.79, 9.37, 12.75)
df <- data.frame(type1, type2)
head(df)
```

| | type1 <dbl> | type2 <dbl> |
|---|----------------|----------------|
| 1 | 3.03 | 3.19 |
| 2 | 5.53 | 4.26 |
| 3 | 5.60 | 4.47 |
| 4 | 9.30 | 4.53 |
| 5 | 9.92 | 4.67 |
| 6 | 12.51 | 4.69 |

6 rows

```
#we reformat the data frame to tidy format for easier analysis
df_tidy<- df %>% gather(key=type, value=time, `type1`, `type2`)
head(df_tidy)
```

| | type <chr> | time <dbl> |
|---|---------------|---------------|
| 1 | type1 | 3.03 |
| 2 | type1 | 5.53 |
| 3 | type1 | 5.60 |
| 4 | type1 | 9.30 |
| 5 | type1 | 9.92 |
| 6 | type1 | 12.51 |

6 rows

✓ b: nonparametric Mann-Whitney test (Wilcoxon rank sum test)

We test the hypothesis with a nonparametric method (Wilcoxin rank sum test).

H_0 : type 1 distribution is the same as type 2 distribution

H_1 : type 1 distribution is different from type 2 distribution.

We compute the rank sum for the smaller size group (in this case both groups are the same) and the smaller rank sum.

```
df_tidy <- df_tidy %>% mutate(ranks=rank(time))
df_tidy
```

| type | time | ranks |
|-------|-------|-------|
| <chr> | <dbl> | <dbl> |
| type1 | 3.03 | 1 |
| type1 | 5.53 | 8 |
| type1 | 5.60 | 9 |
| type1 | 9.30 | 11 |
| type1 | 9.92 | 13 |
| type1 | 12.51 | 14 |
| type1 | 12.95 | 17 |
| type1 | 15.21 | 18 |
| type1 | 16.04 | 19 |
| type1 | 16.84 | 20 |

1-10 of 20 rows

Previous [1](#) [2](#) [Next](#)

```
sum_rank_df <- df_tidy %>% group_by(type) %>% summarize(sum_rank=sum(ranks))
sum_rank_df
```

| type | sum_rank |
|-------|----------|
| <chr> | <dbl> |
| type1 | 130 |
| type2 | 80 |

To test our hypothesis exactly, we use look up table with two sided alpha=.05, $n_1 = n_2 = 10$ on page A22 in Rice. We see that the cutoff is 78. We got 80 so we accept the null.

Lets compare this with the normal approximation of the rank sum. With $m = n = 10$ we have

$$T \sim N(m(m + n + 1)/2, mn(m + n + 1)/12) = N(105, 175).$$

We compute the two sided p-value of rank sum T=80.

```
pval=2*(1-pnorm(abs((80-105))/sqrt(175)))
pval
```

```
## [1] 0.05878172
```

The p-value is slightly larger than the $\alpha = .05$ so we accept the null.

c: t-test for unpaired data

This is 11.6.21a where we do a t-test to compare the means of two normal distributions.

```
#check whether we can assume equal variance  
var(type1)
```

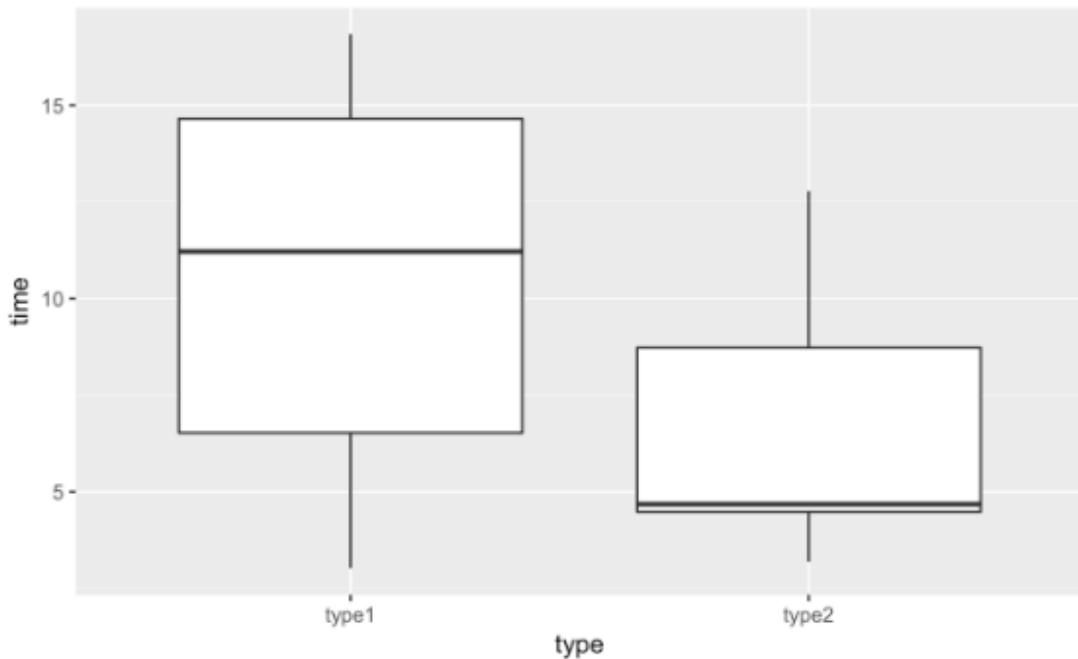
```
## [1] 23.22551
```

```
var(type2)
```

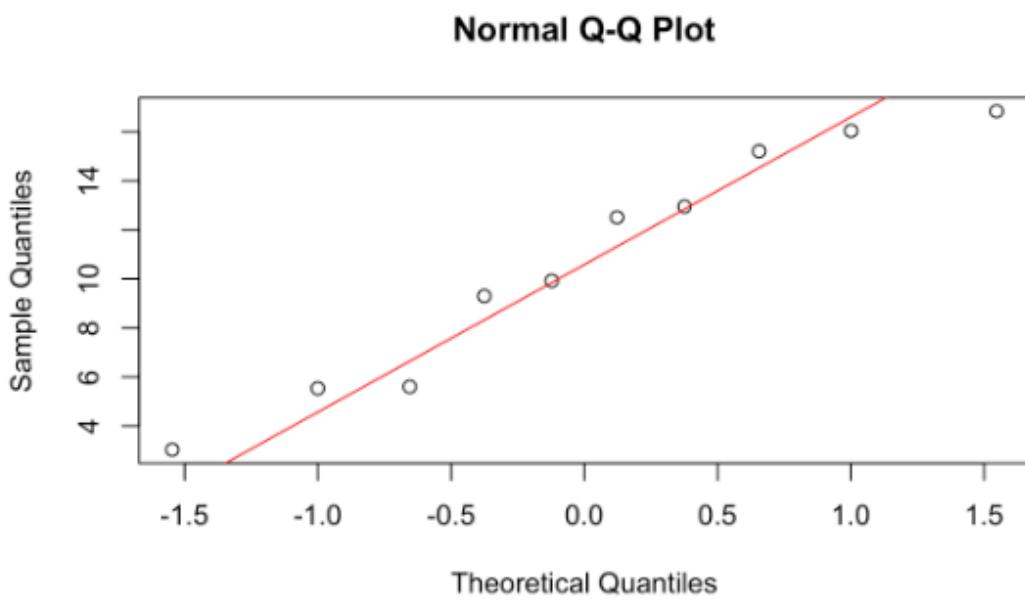
```
## [1] 12.97749
```

Variances appear to be different

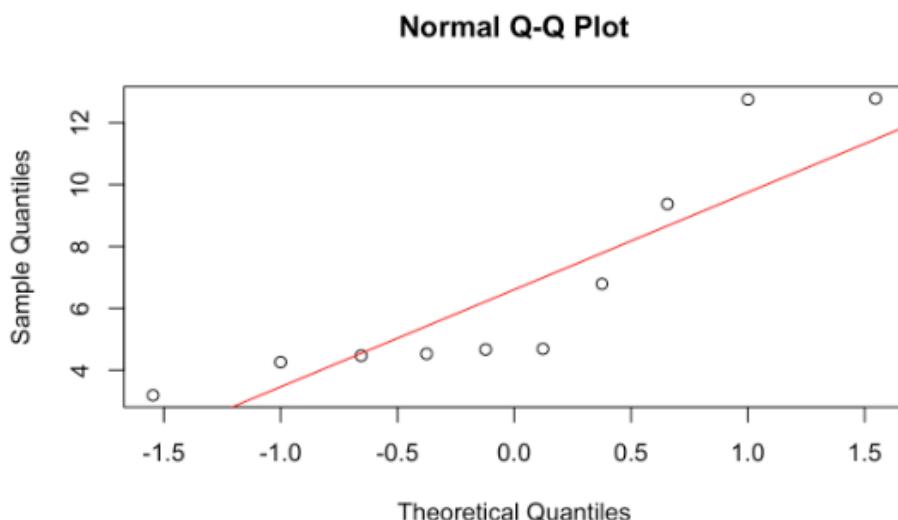
```
df_tidy %>% ggplot(aes(x=type,y=time)) + geom_boxplot()
```



```
qqnorm(type1); qqline(type1,col=2)
```



```
qqnorm(type2); qqline(type2,col=2)
```



```
#default is paired=FALSE, var.equal=FALSE, conf.level=0.95  
t.test(alternative="two.sided",type1,type2)
```

```
##  
## Welch Two Sample t-test  
##  
## data: type1 and type2  
## t = 2.0723, df = 16.665, p-value = 0.05408  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
## -0.07752643 7.96352643  
## sample estimates:  
## mean of x mean of y  
## 10.693 6.750
```

Note that the p-value for the parametric t.test is close to the p-value for the nonparametric Mann-Whitney test. This indicates that the assumptions of the t-test aren't badly violated.

We know that the t-test is uniformly most powerful when the assumptions of the t-test are met (i.e. the rejection region of the t-test isn't a function of the value of the alternative). Under these conditions it is known that the Mann Whitney test is almost as powerful (95% as powerful). When the assumptions of the t-test are not satisfied the results of the t-test can be misleading and we shouldn't really talk about the power of the test. In this case the Mann Whitney test is preferred.

We can also do a Mann-Whitney test with the built in Wilcoxon.test() function in R. Unfortunately this test computes a test statistic, $W = T - m(m + 1)/2 = 80 - 55 = 25$. W is the number of pairs of type1 and type 2 where type2 is less than type1. Out of 100 we would expect 50 if the null is true and we got 25. This isn't too extreme.

```
wilcox.test(type2, type1, alternative="two.sided", exact=FALSE,conf.int=.95, correct=FALSE)
```

```
##  
## Wilcoxon rank sum test  
##  
## data: type2 and type1  
## W = 25, p-value = 0.05878  
## alternative hypothesis: true location shift is not equal to 0  
## 95 percent confidence interval:  
## -8.42001523 0.07004542  
## sample estimates:  
## difference in location  
## -4.27779
```

We accept the null that the two distributions are the same.

Which test is better in this case, t-test or nonparametric test? In my opinion you have a lot to gain for using the Mann-Whitney test and not much to lose.

Sec 11.3.2 Non Parametric Wilcoxon
Signed Rank Test,

For Paired data.

| B | A | D _i | D _i | Rank | Signed Rank |
|----|----|----------------|----------------|------|-------------|
| 25 | 27 | 2 | 2 | 2 | 2 |
| 29 | 25 | -4 | 4 | 3 | -3 |
| 60 | 59 | -1 | 1 | 1 | -1 |
| 27 | 37 | 10 | 10 | 4 | 4 |

$w_+ = 2 + 4 = 6$

Called V in R ↑ sum of pos signed ranks.

H_0 : Distribution of D_i is symmetric around 0

H_1 : " " " " is not symmetric around 0

Then Under the null that D_i are symmetrically distributed around 0

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

Pf / see book,

Here it is in R:

Paired 2 sample tests

This is problem 48 in Rice. Discuss the following code with your neighbor.

The amount of urinary protein before and after a certain treatment is given in the table below:

```
before <- c(24.6,17,16,10.4,8.2,7.9,8.2,7.9,5.8,5.4,5.1,4.7)
after <- c(10.1,5.7,5.6,3.4,6.5,0.7,6.5,0.7,6.1,4.7,2.0,2.9)
diff=before-after
abs_diff=abs(diff)
rank <- rank(abs_diff)
signed_rank <- rank*(diff/abs_diff)
df <- data.frame(before,after,diff, abs_diff,rank, signed_rank)
df
```

| before <dbl> | after <dbl> | diff <dbl> | abs_diff <dbl> | rank <dbl> | signed_rank <dbl> |
|-----------------|----------------|---------------|-------------------|---------------|----------------------|
| 24.6 | 10.1 | 14.5 | 14.5 | 12.0 | 12.0 |
| 17.0 | 5.7 | 11.3 | 11.3 | 11.0 | 11.0 |
| 16.0 | 5.6 | 10.4 | 10.4 | 10.0 | 10.0 |
| 10.4 | 3.4 | 7.0 | 7.0 | 7.0 | 7.0 |
| 8.2 | 6.5 | 1.7 | 1.7 | 3.5 | 3.5 |
| 7.9 | 0.7 | 7.2 | 7.2 | 8.5 | 8.5 |
| 8.2 | 6.5 | 1.7 | 1.7 | 3.5 | 3.5 |
| 7.9 | 0.7 | 7.2 | 7.2 | 8.5 | 8.5 |
| 5.8 | 6.1 | -0.3 | 0.3 | 1.0 | -1.0 |
| 5.4 | 4.7 | 0.7 | 0.7 | 2.0 | 2.0 |

1-10 of 12 rows

Previous [1](#) [2](#) Next

null and alternative:

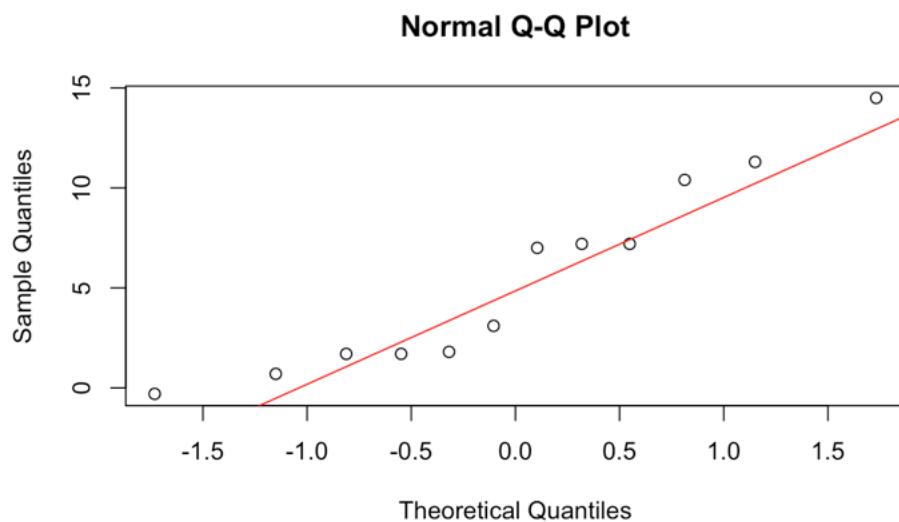
H_0 : the distribution of D_i is symmetric around zero

H_1 : the distribution of D_i isn't symmetric around zero

✓ a) t-test

check assumption that D_i are iid normal:

```
qqnorm(diff); qqline(diff, col=2)
```



If assumptions that the differences looks normal we perform a t-test using command `t.test` and make conclusion:

```
#step2
t.test(before, after, paired = TRUE)
```

```
##
##  Paired t-test
##
## data: before and after
## t = 4.0012, df = 11, p-value = 0.002082
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  2.485818 8.564182
## sample estimates:
## mean of the differences
##                      5.525
```

✓ b) Wilcoxon Signed Rank test

Use the command `wilcox.test()` to perform a Wilcoxon Signed Rank test. Use console below for your answer

Code Start Over Solution Run Code

```
1 wilcox.test(before,after,paired = TRUE,correct=FALSE)
2
3
```

| before <dbl> | after <dbl> | diff <dbl> | abs_diff <dbl> | rank <dbl> | signed_rank <dbl> |
|-----------------|----------------|---------------|-------------------|---------------|----------------------|
| 24.6 | 10.1 | 14.5 | 14.5 | 12.0 | 12.0 |
| 17.0 | 5.7 | 11.3 | 11.3 | 11.0 | 11.0 |
| 16.0 | 5.6 | 10.4 | 10.4 | 10.0 | 10.0 |
| 10.4 | 3.4 | 7.0 | 7.0 | 7.0 | 7.0 |
| 8.2 | 6.5 | 1.7 | 1.7 | 3.5 | 3.5 |
| 7.9 | 0.7 | 7.2 | 7.2 | 8.5 | 8.5 |
| 8.2 | 6.5 | 1.7 | 1.7 | 3.5 | 3.5 |
| 7.9 | 0.7 | 7.2 | 7.2 | 8.5 | 8.5 |
| 5.8 | 6.1 | -0.3 | 0.3 | 1.0 | -1.0 |
| 5.4 | 4.7 | 0.7 | 0.7 | 2.0 | 2.0 |

1-10 of 12 rows

Previous **1** **2** Next

Warning in wilcox.test.default(before, after, paired = TRUE, correct = FALSE): cannot compute exact p-value with ties

```
Wilcoxon signed rank test

data: before and after
V = 77, p-value = 0.002852
alternative hypothesis: true location shift is not equal to 0
```

Note that the test statistic is called **V** in R. Which test is better in this case, t-test or nonparametric test?

If assumptions are met,
t-test is a uniformly most
powerful test. The Wilcoxon
Signed Rank Test is around
95% as powerful.

ex (See 11.6.26b) - HW #8

Let $X_1 \sim N(0, 1)$ } independent samples
 $Y_1, Y_2 \sim N(1, 1)$ }

Determine the variance of the rank sum of the X_i 's.

i.e. $\text{Var}(T_X)$.

$$U_Y = T_Y - \frac{m(m+1)}{2}$$

$$T_X + T_Y = \frac{(m+n)(m+n+1)}{2}$$

$$\text{so } \text{Var}(T_X) = \text{Var}(U_Y)$$

so we find $\text{Var}(U_Y)$.

$$U_Y = \sum_{i=1}^m \sum_{j=1}^n I(X_i < Y_j)$$

$X_1 - Y_1, X_2 - Y_2$
are dependent

$$\text{Var}(U_Y) = \text{Var}(\sum I(X_i < Y_1) + I(X_i < Y_2))$$

$$= \text{Var}(I(X_i < Y_1)) + \text{Var}(I(X_i < Y_2)) + \text{Cov}(I(X_i < Y_1), I(X_i < Y_2))$$

$$\text{Cov}(I(X_i < Y_1), I(X_i < Y_2)) = E(I(X_i < Y_1)I(X_i < Y_2)) - E(I(X_i < Y_1))E(I(X_i < Y_2))$$

$$E(I(X_i < Y_1)I(X_i < Y_2)) = P(X_i < Y_1 \text{ and } X_i < Y_2)$$

hard to find analytically
so let's approximate it by
simulation in R.

R Simulations

```
B=10000

fun <- function(){
  a <- rnorm(1)
  b <- rnorm(2, mean = 1)
  p <- prod(a < b) — I(X_i < Y_1).I(X_i < Y_2)
}

vec <- replicate(B, fun())
mean(vec)
```

You will find that $P(X_i < Y_1 \text{ and } X_i < Y_2) = .634$

$$\text{and } \text{Cov}(I(X_i < Y_1), I(X_i < Y_2)) = .634 - (.76)^2 \approx .056$$

$$\boxed{\text{Var}(U_Y) = 2(.76)(.24) - 2(.056)}$$

Last time

Sec 11.2 and 11.3

nonparametric tests for
comparing the distribution of
2 continuous populations

Today

Sec 12.2 Analysis of Variance (ANOVA)

parametric test (Normal theory
F-test) comparing the means of
2 or more continuous populations.

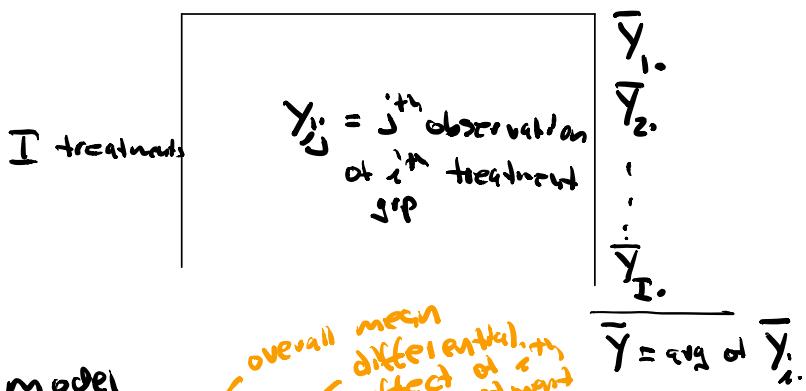
t.

Sec 12.2 ANOVA (Analysis of Variance)

One-way layout is an experimental design in which independent measurements are made under each of several treatments.

Picture

J observations



model

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \sim N(\mu + \alpha_i, \sigma^2)$$

overall mean
 differential effect of ith treatment

$$\sum_{i=1}^I \alpha_i = 0$$

$$Y_{ij} \sim N(\mu + \alpha_i, \sigma^2)$$

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_I = 0 \quad E(Y_{ij}) = \mu$$

$$H_1: \text{not all } \alpha_i \text{ are zero.}$$

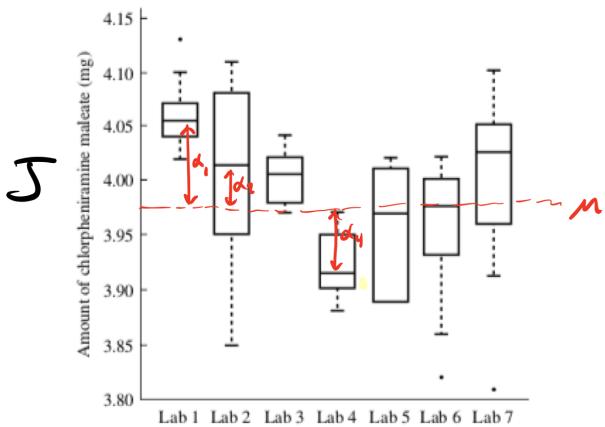


FIGURE 12.1 Boxplots of determinations of amounts of chlorpheniramine in tablets by seven laboratories.

A nova

Test null that differential effect $\alpha_i \rightarrow$ zero,

ANOVA is based on decomposition of the total variance (SST) (sum of squares total) of your data into group variance (SSW) (sum of squares within) and the intergroup variance (SSB) (sum of squares between),

$$SST = SSW + SSB$$

$$\sum_{i=1}^J \sum_{j=1}^n (Y_{ij} - \bar{Y})^2 \quad \sum_{i=1}^J \sum_{j=1}^n (Y_{ij} - \bar{Y}_{i.})^2 \quad \sum_{i=1}^J \sum_{j=1}^n (\bar{Y}_{i.} - \bar{Y})^2$$

The algebraic proof is on p408,

SSW and SSB are independent RVs,
to see this:

Facts from Chap 6

$$x_1, \dots, x_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

① \bar{X} and $x_j - \bar{X}$ are indep RVs
(thm A P195)

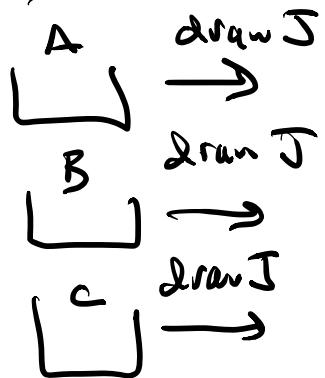
② $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$ (thm B P197)

$SSB = \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij} - \bar{Y})^2$ is a function of \bar{Y}_{ij} since $\bar{Y} = \frac{1}{I} \sum_{i=1}^I \bar{Y}_{ij}$

$SSW = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{Y}_{ij})^2$ is a function of $y_{ij} - \bar{Y}_{ij}$

$\Rightarrow SSW$ and SSB are indep RVs.

ex $I=3, J=3$



$J=3$

| | | | $\bar{Y}_{1.} = 2$ |
|---|---|---|--------------------|
| | | | $\bar{Y}_{2.} = 4$ |
| | | | $\bar{Y}_{3.} = 6$ |
| A | 3 | 2 | 1 |
| B | 5 | 3 | 4 |
| C | 5 | 6 | 7 |

$\bar{Y} = 4$

$$SST = \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y})^2$$

$$= (3-4)^2 + (2-4)^2 + (1-4)^2$$

$$+ (5-4)^2 + (3-4)^2 + (4-4)^2$$

$$+ (5-4)^2 + (6-4)^2 + (7-4)^2$$

$$= 30$$

this is where
we assume
the null.

Under the null we have $I \cdot J$

draws from a single box,

$$\frac{SST}{IJ-1} = S^2 \Rightarrow \frac{(IJ-1)S^2}{\sigma^2} \sim \chi^2_{IJ-1}$$

We say SST has $IJ-1$ d.f

assuming the null

SSW

$J=3$

| | | | | |
|---|---|---|---|--------------------|
| A | 3 | 2 | 1 | $\bar{Y}_{1.} = 2$ |
| B | 5 | 3 | 4 | $\bar{Y}_{2.} = 4$ |
| C | 5 | 6 | 7 | $\bar{Y}_{3.} = 6$ |

$\bar{Y} = 4$

$$\begin{aligned}
 SSW &= \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y}_{i.})^2 \\
 &= (3-2)^2 + (2-2)^2 + (1-2)^2 \\
 &\quad + (5-4)^2 + (3-4)^2 + (4-4)^2 \\
 &\quad + (5-6)^2 + (3-6)^2 + (7-6)^2 \\
 &= 6 \leftarrow (J-1)S_A^2 + (J-1)S_B^2 + (J-1)S_C^2 \\
 \frac{SSW}{\sigma^2} &= I(J-1) \frac{(J-1)S_A^2 + (J-1)S_B^2 + (J-1)S_C^2}{I(J-1)} = S_{\text{pooled}}^2 \\
 &\sim \chi_{I(J-1)}^2
 \end{aligned}$$

We say SSW has $I(J-1)$ d.f.

SSB

$\bar{J} = 3$

| | | | | |
|---|---|---|---|--------------------|
| A | 3 | 2 | 1 | $\bar{Y}_{1.} = 2$ |
| B | 5 | 3 | 4 | $\bar{Y}_{2.} = 4$ |
| C | 5 | 6 | 7 | $\bar{Y}_{3.} = 6$ |

$\bar{Y} = 4$

$$\begin{aligned}
 SSB &= \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij} - \bar{Y})^2 \\
 &= (2-4)^2 + (2-4)^2 + (2-4)^2 \\
 &\quad + (4-4)^2 + (4-4)^2 + (4-4)^2 \\
 &\quad + (6-4)^2 + (6-4)^2 + (6-4)^2 = \boxed{24}
 \end{aligned}$$

Note:

$$\begin{array}{ccc}
 SSJ & = & SSW + SSB \\
 \parallel & & \parallel \\
 30 & & 24
 \end{array}$$

✓

Facts

1) $SST = SSW + SSB$

2) SSW and SSB are indep

3) Chi square subtraction property

If $X = X_1 + X_2$ with X_1, X_2 indep

and $X \sim \chi^2_{n_1+n_2}$

$X_1 \sim \chi^2_{n_1}$

then $X_2 \sim \chi^2_{n_2}$.

here $n_1 = \text{df of } SSW = I(J-1)$

$n_1+n_2 = \text{df of } SST = IJ-1$

$$\Rightarrow n_2 = \text{df of } SSB = IJ-1 - I(J-1) \\ = \boxed{I-1}$$

We say SSB has $I-1$ d.f. assuming null.

(we need null assumption since we use fact that SST has $IJ-1$ d.f. under null.).

F distribution — Chap 6 pg 94

Recall χ^2 , t, F are distributions based on the normal distributions.

Defⁿ Let U and V be indep χ^2 RV w/ a and b degrees of freedom, the distribution,

$$F_{a,b} = \frac{U/a}{V/b} \rightarrow \text{called the } \boxed{\text{F distribution}}$$

F test (assuming H₀ null) $\frac{SSB/\sigma^2}{SSW/\sigma^2}$

$$F_{I-1, I(J-1)} = \frac{\frac{SSB/\sigma^2}{I-1}}{\frac{SSW/\sigma^2}{I(J-1)}}$$

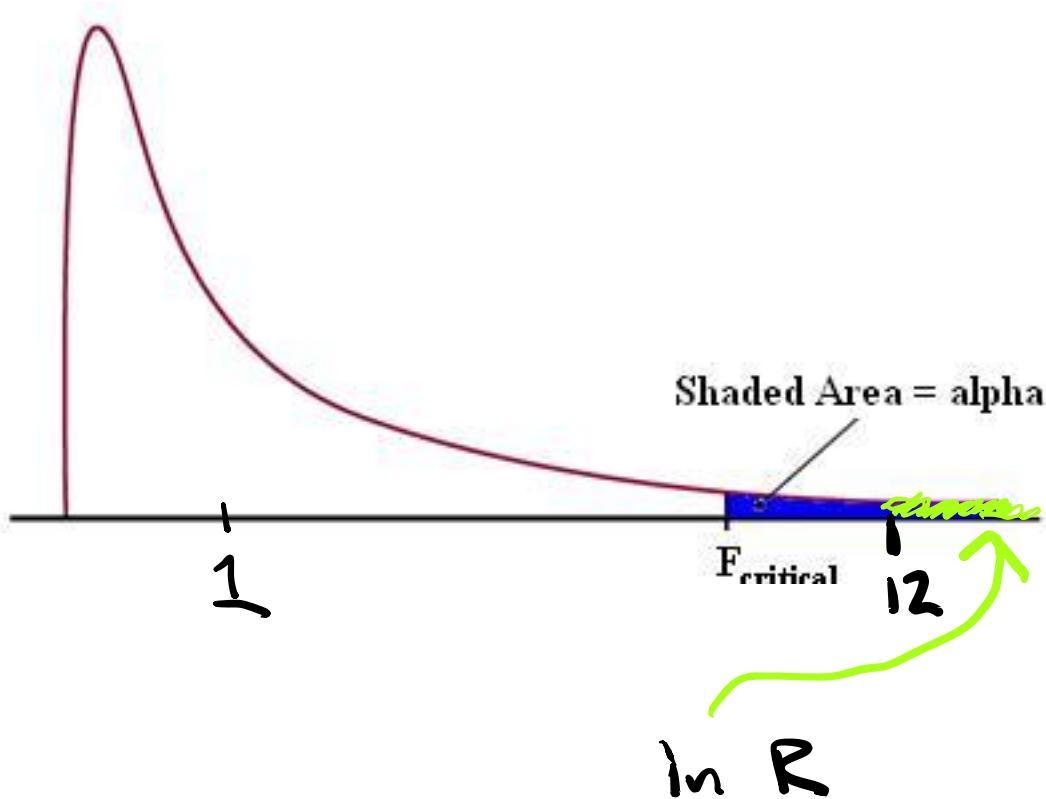
↑ between d.f. ↑ within d.f.

$$= \frac{\cancel{SSB/I-1}}{\cancel{SSW/I(J-1)}}$$

For UV example

$$F_{2,6} = \frac{27/2}{6/6} = 12$$

$F_{2,6}$ distribution.



$$1 - Pf(12, 2, 6) = .008$$

↑ K df within
df betw

so we reject null
that all treatments have same
mean.

Stat 135 Lec 30

Last time

Sec 12.2 One way ANOVA with different size grps J_1, J_2, \dots, J_I

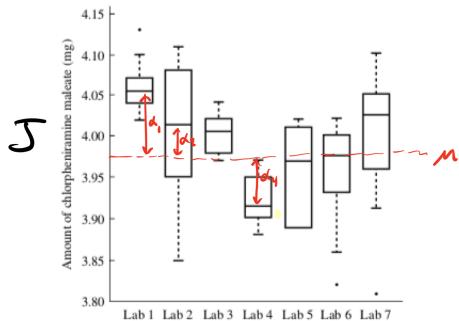


FIGURE 12.1 Boxplots of determinations of amounts of chlorpheniramine in tablets by seven laboratories.

$$\text{model } Y_{ij} = \mu + \alpha_i + \epsilon_{ij} \sim N(0, \sigma^2)$$

overall mean
 differential, i, effect of treatment

A nova

$$SST = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y})^2 \quad \text{and} \quad \frac{SST}{\sigma^2} \sim \chi^2_{\sum_{i=1}^I (J_i - 1)}$$

assuming null
that all $\alpha_i = 0$

$$SSW = \sum_{i=1}^I \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{i\cdot})^2 \quad \text{and} \quad \frac{SSW}{\sigma^2} \sim \chi^2_{\sum_{i=1}^I (J_i - 1)}$$

$$SSB = \sum_{i=1}^I \sum_{j=1}^{J_i} (\bar{Y}_{i\cdot} - \bar{Y})^2 \quad \text{and} \quad \frac{SSB}{\sigma^2} \sim \chi^2_{I-1} \quad \text{assuming true null}$$

$$\frac{F_{\text{test}}}{F_{(I-1, \sum_{i=1}^I (J_i - 1))}} = \frac{\frac{SSB}{I-1}}{\frac{SSW}{\sum_{i=1}^I (J_i - 1)}} \quad \text{assuming true null}$$

① Assumptions

- Today
- ① Shortcut Bonferroni
 - ② Kruskal Wallis (KW) test (non parametric)
 - ③ KW test in R,

① Assumptions

The normal theory F-test has several assumptions:

- 1) Our data y_{ij} is normal
- 2) The variance of each treatment group is the same σ^2
- 3) All observations are independent.

If we don't satisfy these use non-parametric test (Kruskall-Wallis test).

① Shortcut Bonferroni test.

Bonferroni test

Suppose the F test rejects the null (i.e. at least one pair of means are different)

If we do $\binom{I}{2}$ simultaneous pairwise t-tests to see which pairs of treatment groups have diff. means, the type I error $P_0(\delta(x)=1)$ is additive and it is quite likely that we reject the null for one of the t-tests.

Bonferroni's method is to make significance level of each test $\frac{\alpha}{\binom{I}{2}}$ level of significance.

Shortcut Bonferroni method

The Bonferroni method is a lot of work and the small significance level makes it hard to reject the null for each t-test.

Before I do the Bonferroni test I do this shortcut test:

For each treatment group $H_0: \mu_i = 0$ ($\mu_i = \mu$)
 $H_1: \mu_i \neq 0$

We compute 3 simultaneous $100(1 - \frac{\alpha}{3})$ two-sided CI for the true mean μ_i for each group.

We simply look whether one of the CI doesn't overlap with the others.

On test or HW I will accept this method.

Ex In our 3 grp example

| | | | | | |
|---|---|---|---|-----------------|--|
| A | 3 | 2 | 1 | $\bar{Y}_1 = 2$ | $(J-1)S_A^2 = (3-2)^2 + (2-2)^2 + (1-2)^2 = 2$ |
| B | 5 | 3 | 4 | $\bar{Y}_2 = 4$ | $(J-1)S_B^2 = 2$ |
| C | 5 | 6 | 7 | $\bar{Y}_3 = 6$ | $(J-1)S_C^2 = 2$ |
| | | | | $\bar{Y} = 4$ | |

$$F_{2,6} = 1008 < 105 \Rightarrow \text{reject null}$$

$"\alpha"$

So we perform a Shortcut Bonferroni.

Recall

$$\frac{\bar{Y}_{i.} - M_i}{S/\sqrt{J_i}} \sim t_{J_i-1}$$

\Rightarrow A $100(1 - \frac{\alpha/2}{I})\%$ CI for M_i

$$\therefore \bar{Y}_{i.} \pm t_{J_i-1} \left(\frac{\alpha/2}{I} \right) \cdot S/\sqrt{J_i}$$

In our example we had $S_i = 1$.

For $\frac{\alpha}{I} = .025$ there a special

Bonferroni table $t_{J_i-1} \left(\frac{\alpha/2}{I} \right)$

so $\bar{Y}_1 \pm t_2 \left(\frac{.025}{3} \right) \cdot \frac{1}{\sqrt{3}} = 2 \pm 4.41 = 7.648 \left(\frac{1}{\sqrt{3}} \right)$

in table next page

$$\bar{Y}_2 \pm t_2 \left(\frac{.025}{3} \right) \cdot \frac{1}{\sqrt{3}} = 4 \pm 4.41$$

$$\bar{Y}_3 \pm t_2 \left(\frac{.025}{3} \right) \cdot \frac{1}{\sqrt{3}} = 6 \pm 4.41$$

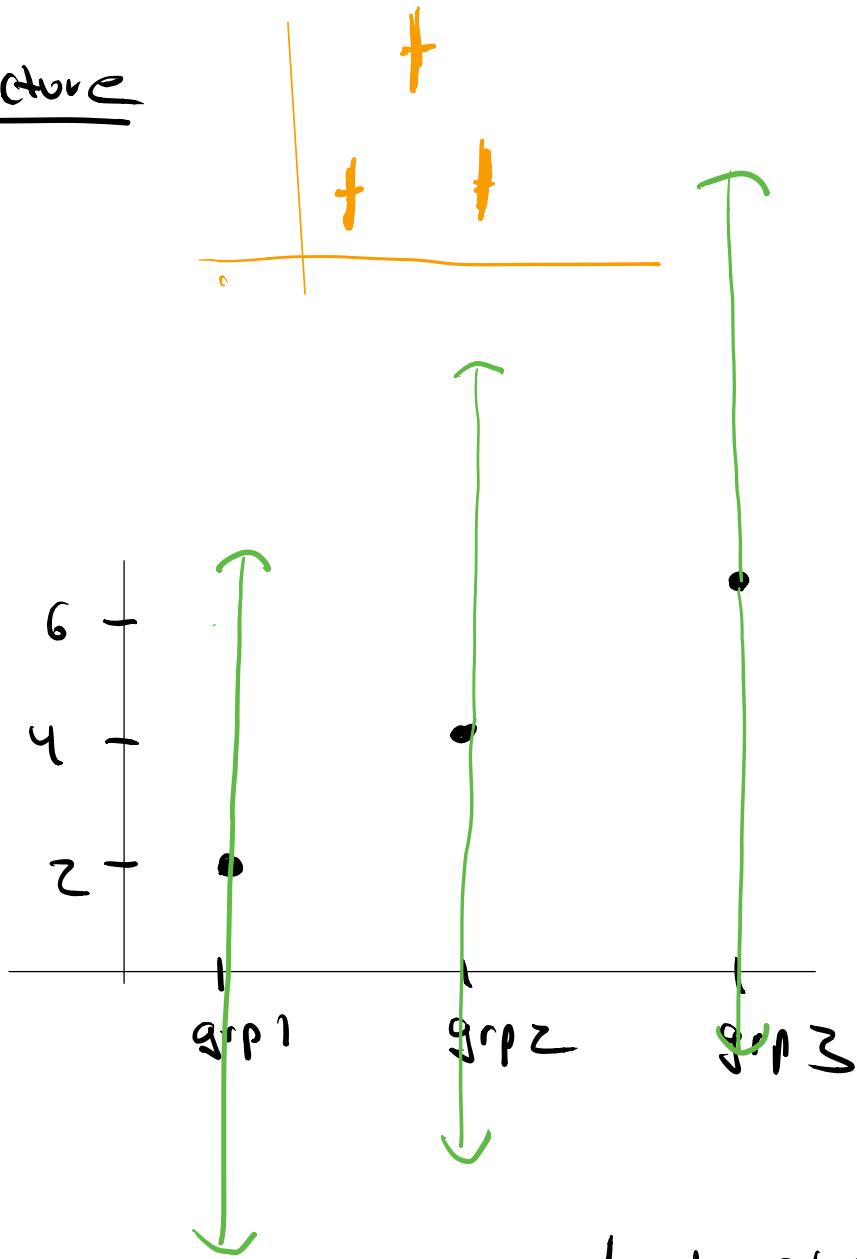
Bonferroni table (not in book)

TABLE 11 Bonferroni Multipliers for 95% Confidence Intervals

The values given in the table are $t_{df, 0.025/k}$ where k is the number of tests.

| df | NUMBER OF TESTS | | | | | | | | | |
|----------|-----------------|--------|--------|--------|--------|--------|---------|---------|---------|---------|
| | 1 | 2 | 3 | 4 | 5 | 6 | 8 | 10 | 15 | 20 |
| 1 | 12.706 | 25.452 | 38.185 | 50.923 | 63.657 | 76.384 | 101.856 | 127.321 | 190.946 | 254.647 |
| 2 | 4.303 | 6.205 | 7.648 | 8.860 | 9.925 | 10.885 | 12.590 | 14.089 | 17.275 | 19.963 |
| 3 | 3.182 | 4.177 | 4.857 | 5.392 | 5.841 | 6.231 | 6.895 | 7.453 | 8.575 | 9.465 |
| 4 | 2.776 | 3.495 | 3.961 | 4.315 | 4.604 | 4.851 | 5.261 | 5.598 | 6.254 | 6.758 |
| 5 | 2.571 | 3.163 | 3.534 | 3.810 | 4.032 | 4.219 | 4.526 | 4.773 | 5.247 | 5.604 |
| 6 | 2.447 | 2.969 | 3.287 | 3.521 | 3.707 | 3.863 | 4.115 | 4.317 | 4.698 | 4.981 |
| 7 | 2.365 | 2.841 | 3.128 | 3.335 | 3.499 | 3.636 | 3.855 | 4.029 | 4.355 | 4.595 |
| 8 | 2.306 | 2.752 | 3.016 | 3.206 | 3.355 | 3.479 | 3.677 | 3.833 | 4.122 | 4.334 |
| 9 | 2.262 | 2.685 | 2.933 | 3.111 | 3.250 | 3.364 | 3.547 | 3.690 | 3.954 | 4.146 |
| 10 | 2.228 | 2.634 | 2.870 | 3.038 | 3.169 | 3.277 | 3.448 | 3.581 | 3.827 | 4.005 |
| 11 | 2.201 | 2.593 | 2.820 | 2.981 | 3.106 | 3.208 | 3.370 | 3.497 | 3.728 | 3.895 |
| 12 | 2.179 | 2.560 | 2.779 | 2.934 | 3.055 | 3.153 | 3.308 | 3.428 | 3.649 | 3.807 |
| 13 | 2.160 | 2.533 | 2.746 | 2.896 | 3.012 | 3.107 | 3.256 | 3.372 | 3.584 | 3.735 |
| 14 | 2.145 | 2.510 | 2.718 | 2.864 | 2.977 | 3.069 | 3.214 | 3.326 | 3.529 | 3.675 |
| 15 | 2.131 | 2.490 | 2.694 | 2.837 | 2.947 | 3.036 | 3.177 | 3.286 | 3.484 | 3.624 |
| 16 | 2.120 | 2.473 | 2.673 | 2.813 | 2.921 | 3.008 | 3.146 | 3.252 | 3.444 | 3.581 |
| 17 | 2.110 | 2.458 | 2.655 | 2.793 | 2.898 | 2.984 | 3.119 | 3.222 | 3.410 | 3.543 |
| 18 | 2.101 | 2.445 | 2.639 | 2.775 | 2.878 | 2.963 | 3.095 | 3.197 | 3.380 | 3.510 |
| 19 | 2.093 | 2.433 | 2.625 | 2.759 | 2.861 | 2.944 | 3.074 | 3.174 | 3.354 | 3.481 |
| 20 | 2.086 | 2.423 | 2.613 | 2.744 | 2.845 | 2.927 | 3.055 | 3.153 | 3.331 | 3.455 |
| 25 | 2.060 | 2.385 | 2.566 | 2.692 | 2.787 | 2.865 | 2.986 | 3.078 | 3.244 | 3.361 |
| 30 | 2.042 | 2.360 | 2.536 | 2.657 | 2.750 | 2.825 | 2.941 | 3.030 | 3.189 | 3.300 |
| 40 | 2.021 | 2.329 | 2.499 | 2.616 | 2.704 | 2.776 | 2.887 | 2.971 | 3.122 | 3.227 |
| 50 | 2.009 | 2.311 | 2.477 | 2.591 | 2.678 | 2.747 | 2.855 | 2.937 | 3.083 | 3.184 |
| 60 | 2.000 | 2.299 | 2.463 | 2.575 | 2.660 | 2.729 | 2.834 | 2.915 | 3.057 | 3.156 |
| 70 | 1.994 | 2.291 | 2.453 | 2.564 | 2.648 | 2.715 | 2.820 | 2.899 | 3.039 | 3.137 |
| 80 | 1.990 | 2.284 | 2.445 | 2.555 | 2.639 | 2.705 | 2.809 | 2.887 | 3.026 | 3.122 |
| 100 | 1.984 | 2.276 | 2.435 | 2.544 | 2.626 | 2.692 | 2.793 | 2.871 | 3.007 | 3.102 |
| 140 | 1.977 | 2.266 | 2.423 | 2.530 | 2.611 | 2.676 | 2.776 | 2.852 | 2.986 | 3.079 |
| 1000 | 1.962 | 2.245 | 2.398 | 2.502 | 2.581 | 2.643 | 2.740 | 2.813 | 2.942 | 3.031 |
| ∞ | 1.960 | 2.241 | 2.394 | 2.498 | 2.576 | 2.638 | 2.734 | 2.807 | 2.935 | 3.023 |

Picture



No conclusion. The $\frac{d}{I}$ level of significance is too big. We need I to be bigger to see non overlapping intervals.

Stat135 lecture 31

(data cleaning)

a: F-test

b: nonparametric Kruskal-Wallis test

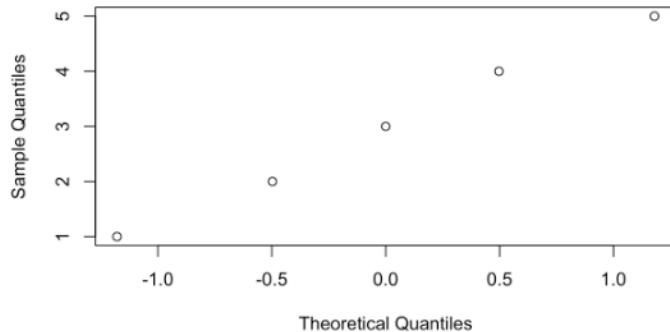
```
library(tidy)
#time to failure for three types of engine bearings

treatment1 <- c(1,2,3,4,5)
treatment2 <- c(6,7,8,9)
treatment3 <- c(10,11,12)

#you should check that normality assumption is satisfied.
qgnorm(treatment1)
```

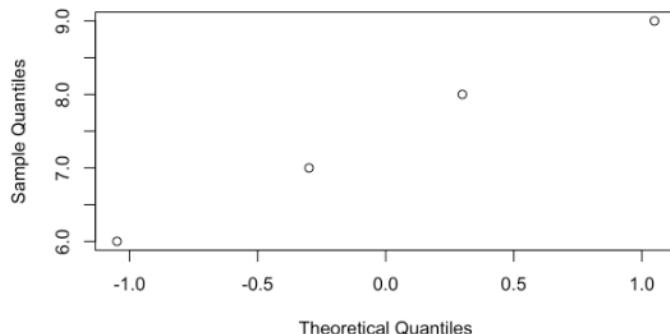
Start Over

Normal Q-Q Plot



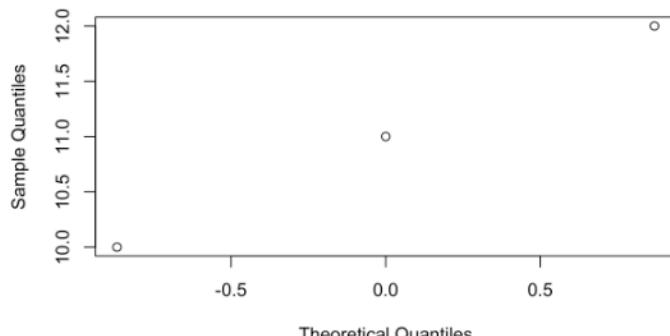
```
qgnorm(treatment2)
```

Normal Q-Q Plot



```
qgnorm(treatment3)
```

Normal Q-Q Plot

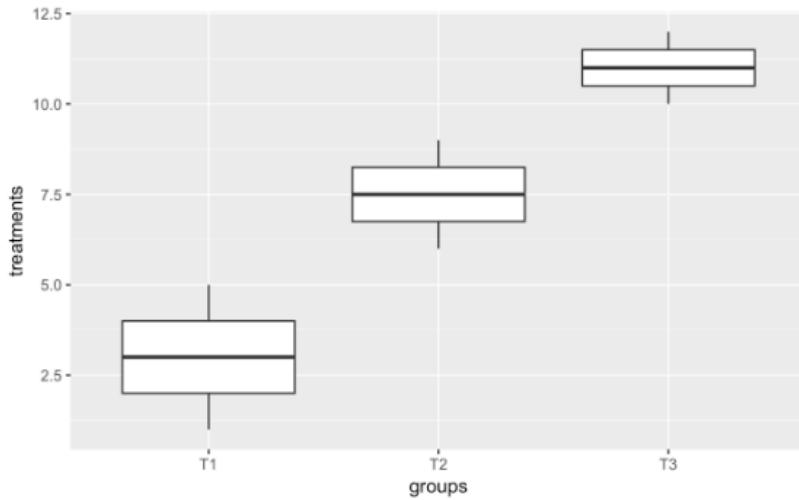


```
groups <- c(rep("T1",5),rep("T2",4),rep("T3",3))
treatments <- c(treatment1,treatment2,treatment3)
df_tidy <- data.frame(groups,treatments)
df_tidy
```

| groups | treatments |
|-----------------|---|
| <fctr> | <dbl> |
| T1 | 1 |
| T1 | 2 |
| T1 | 3 |
| T1 | 4 |
| T1 | 5 |
| T2 | 6 |
| T2 | 7 |
| T2 | 8 |
| T2 | 9 |
| T3 | 10 |
| 1-10 of 12 rows | Previous 1 2 Next |

You should make boxplot to examine the means and variances

```
df_tidy %>% ggplot(aes(x=groups,y=treatments))+ geom_boxplot()
```



a: F-test

```
I=3
J=c(5,4,3)

mean_T1 <- mean(treatment1)
mean_T2 <- mean(treatment2)
mean_T3 <- mean(treatment3)
mean_T1

## [1] 3

mean_T2

## [1] 7.5

mean_T3

## [1] 11

mean_tot <- mean(df_tidy$treatments)

var_T1 <- var(treatment1)
var_T2 <- var(treatment2)
var_T3 <- var(treatment3)
var_T1

## [1] 2.5

var_T2

## [1] 1.666667

var_T3

## [1] 1

SSW <- (J[1]-1)*var(treatment1) + (J[2]-1)*var(treatment2) + (J[3]-1)*var(treatment3)
SSW

## [1] 17

SSB <- J[1]*(mean_T1-mean_tot)^2 + J[2]*(mean_T2-mean_tot)^2 + J[3]*(mean_T3-mean_tot)^2
SSB

## [1] 126

F <- (SSB/2)/(SSW/9)
F

## [1] 33.35294

1-pf(F,2,9)

## [1] 6.886649e-05
```

Using `oneway.test()`

```
oneway.test(treatments~groups, var.equal = TRUE) #in formula rhs: treatment values lhs: group type  
  
##  
## One-way analysis of means  
##  
## data: treatments and groups  
## F = 33.353, num df = 2, denom df = 9, p-value = 6.887e-05
```

Or can assume not equal variance:

```
oneway.test(treatments~groups)  
  
##  
## One-way analysis of means (not assuming equal variances)  
##  
## data: treatments and groups  
## F = 34.489, num df = 2.000, denom df = 5.771, p-value = 0.000617
```

Shortcut Bonferroni method

Since we reject the null lets do the shortcut Bonferonni method (level of significance for each test is 91.7%) to see which group means are significantly different.

```
#T1  
mean_T1=qt(1-.025/3,4)*sqrt(var_T1)/sqrt(5)
```

```
## [1] 0.199301
```

```
mean_T1+qt(1-.025/3,4)*sqrt(var_T1)/sqrt(5)
```

```
## [1] 5.800699
```

T_1 CI is (0.2,5.8)

```
#T2  
mean_T2=qt(1-.025/3,3)*sqrt(var_T2)/sqrt(4)
```

```
## [1] 4.365041
```

```
mean_T2 +qt(1-.025/3,3)*sqrt(var_T2)/sqrt(4)
```

```
## [1] 10.63496
```

T_2 CI is (4.4, 10.6)

```
#T3  
mean_T3=qt(1-.025/3,2)*sqrt(var_T3)/sqrt(3)
```

```
## [1] 6.583961
```

```
mean_T3 +qt(1-.025/3,2)*sqrt(var_T3)/sqrt(3)
```

```
## [1] 15.41604
```

T_3 CI is (6.6., 15.4).

We see T1 and T3 CIs don't overlap. We conclude that the means of T_1 and T_3 are different.

Lec 3)

Quiz 4 Friday Nov 22

main Whitney Test (not the ANOVA/Bonferroni)
Kruskal Wallis test
Signed rank sum test for Paired data

Last time One way ANOVA F-test
and Bonferroni's shortest comparison test

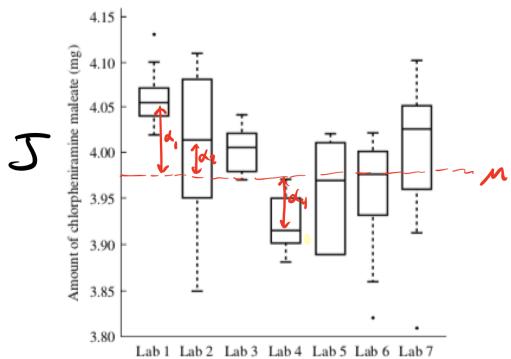


FIGURE 12.1 Boxplots of determinations of amounts of chlorpheniramine in tablets by seven laboratories.

A n o v a

Assumptions of F-test

- 1) Our data y_{ij} is normal
- 2) The variance of each treatment group is the same σ^2
- 3) All observations are independent,

model

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij} \sim N(0, \sigma^2)$$

overall mean differential effect of treatment

$$SSB = \sum_{j=1}^{J_0} \sum_{i=1}^{J_j} (\bar{Y}_{ij} - \bar{Y})^2 \quad \text{and} \quad \frac{SSB}{\sigma^2} \sim \chi^2_{J_0 - 1}$$

Today ① Nonparametric Kruskall Wallis test sec 12.2.3

② Matrix notation in multiple regression

③ Simple standard regression model ($k=1$)

④ Formula for regression coefficients

} sec 14.2

Sec 12.2.3 Kruskal Wallis (KW) test

$$R_{ij} \in \{1, 2, \dots, N\}$$

R_{ij} = rank of y_{ij} in the combined treatment grps.

Null: each treatment grp has same distribution
 Alt: not all treatment grp have same dist.

let $SSB = \sum_{i=1}^I J_i (\bar{R}_{i\cdot} - \bar{R})^2$ be the between grp sum of squares

If the null is true that each treatment grp has the same distribution then SSB should be small,

As before $\frac{SSB}{\sigma^2} \sim \chi^2_{I-1}$

σ^2 can show that
 $\approx \frac{N(Nm)}{12}$

```
library(tidyverse)
#time to failure for three types of engine bearings

treatment1 <- c(1,2,3,4,5)
treatment2 <- c(6,7,8,9)
treatment3 <- c(10,11,12)

#you should check that normality assumption is satisfied.
qqnorm(treatment1)
```

```
library(dplyr)
library(ggplot2)
```

(data cleaning)

a: F-test

b: nonparametric Kruskal-Wallis test

Start Over

b: nonparametric Kruskal-Wallis test

We can decide whether the treatments have an effect without assuming they are normal and have equal variances.

We test the hypothesis with a nonparametric method (Kruskal-Wallis test) which is similar to the Mann-Whitney rank sum test.

H_0 : All the treatment distributions are the same.

H_1 : At least one of the treatment distributions is different.

Step 1

We rank each observation without regard for treatment group, beginning with a rank 1 for the smallest observation. We compute the mean rank for each treatment and the mean of those means. Finally we square the difference of each treatment's average with the total mean.

```
df_tidy <- df_tidy %>% mutate(ranks=rank(treatments))
df_tidy
```

| groups | treatments | ranks |
|--------|------------|-------|
| <fctr> | <dbl> | <dbl> |
| T1 | 1 | 1 |
| T1 | 2 | 2 |
| T1 | 3 | 3 |
| T1 | 4 | 4 |
| T1 | 5 | 5 |
| T2 | 6 | 6 |
| T2 | 7 | 7 |
| T2 | 8 | 8 |
| T2 | 9 | 9 |
| T3 | 10 | 10 |

1-10 of 12 rows

Previous 1 2 Next

```
mean_rank_df <- df_tidy %>% group_by(groups) %>% summarize(mean_rank=mean(ranks))
mean_rank_df
```

| groups | mean_rank |
|--------|-----------|
| <fctr> | <dbl> |
| T1 | 3.0 |
| T2 | 7.5 |
| T3 | 11.0 |

3 rows

```
joined_df <- df_tidy %>%
  left_join(mean_rank_df) %>%
  mutate(tot_mean=mean(mean_rank)) %>%
  mutate(diff_sq=(mean_rank-tot_mean)^2)
joined_df
```

| groups | treatments | ranks | mean_rank | tot_mean | diff_sq |
|--------|------------|-------|-----------|----------|---------|
| <fctr> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| T1 | 1 | 1 | 3.0 | 6.5 | 12.25 |
| T1 | 2 | 2 | 3.0 | 6.5 | 12.25 |
| T1 | 3 | 3 | 3.0 | 6.5 | 12.25 |
| T1 | 4 | 4 | 3.0 | 6.5 | 12.25 |
| T1 | 5 | 5 | 3.0 | 6.5 | 12.25 |
| T2 | 6 | 6 | 7.5 | 6.5 | 1.00 |
| T2 | 7 | 7 | 7.5 | 6.5 | 1.00 |
| T2 | 8 | 8 | 7.5 | 6.5 | 1.00 |
| T2 | 9 | 9 | 7.5 | 6.5 | 1.00 |
| T3 | 10 | 10 | 11.0 | 6.5 | 20.25 |

1-10 of 12 rows

Previous 1 2 Next

Step 2

We compute the Kruskal-Wallis test statistic, K , to obtain a normalized measure of how much the average ranks within each treatment group deviate from the average rank of all the observations.

```
N <- sum(J)
N

## [1] 12

SSB <- sum(joined_df$diff_sq) #sum of square between
SSB

## [1] 126

K <- (12/(N*(N+1)))*sum(joined_df$diff_sq)
K

## [1] 9.692308

1-pchisq(K,df=2)

## [1] 0.007858545
```

Using `Kruskal.test()`

We can also do a Kruskal-Wallis test with the built in `kruskal.test()` function in R.

```
treatments

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

groups <- c(rep(1,5),rep(2,4),rep(3,3))
groups

## [1] 1 1 1 1 1 2 2 2 2 3 3 3

kruskal.test(treatments~groups)

##
##  Kruskal-Wallis rank sum test
##
##  data: treatments by groups
##  Kruskal-Wallis chi-squared = 9.6923, df = 2, p-value = 0.007859
```

We ~~reject~~ accept the null that the ~~two~~ three distributions are the same.

(2)

consider matrices:

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad x_0 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \quad x_1 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{n1} \end{bmatrix}_{n \times 1} \quad \dots \quad x_k = \begin{bmatrix} x_{1k} \\ \vdots \\ x_{nk} \end{bmatrix}_{n \times 1} \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix}_{(k+1) \times 1} \quad e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nk} \end{bmatrix}_{n \times (k+1)} \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix}_{(k+1) \times 1} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \dots + \beta_k x_{1k} + e_1 \\ \vdots \\ \beta_0 + \beta_1 x_{n1} + \dots + \beta_k x_{nk} + e_n \end{bmatrix}_{n \times 1}$$

$$X = \begin{bmatrix} x_0, \dots, x_k \end{bmatrix}_{n \times (k+1)}$$

design matrix

 $n \geq k+1$

X full rank

$\hookrightarrow x_0, x_1, \dots, x_k$ column vectors are independent

so $X'X$ is $(k+1) \times (k+1)$ invertible matrix

The standard statistical model is

$$y = X\beta + e \quad \text{The } e_i \text{ are independent RVs}$$

$\hookrightarrow E(e_i) = 0, \text{Var}(e_i) = \sigma^2$

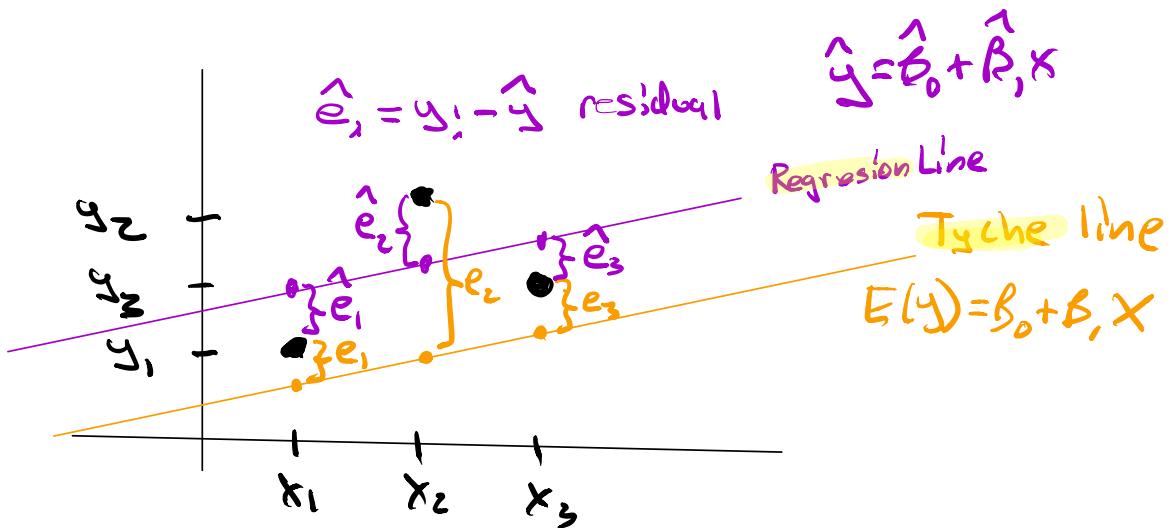
We say e is homoscedastic since $\text{Var}(e_i)$ doesn't depend on X .

X_i are fixed vectors

ER (Simple standard statistical model)

$K=1$
 $n=3$ Let's assume there is a linear relationship between $x = \text{father's height}$ and $y = \text{son's height}$

$$y = \beta_0 + \beta_1 x + e$$

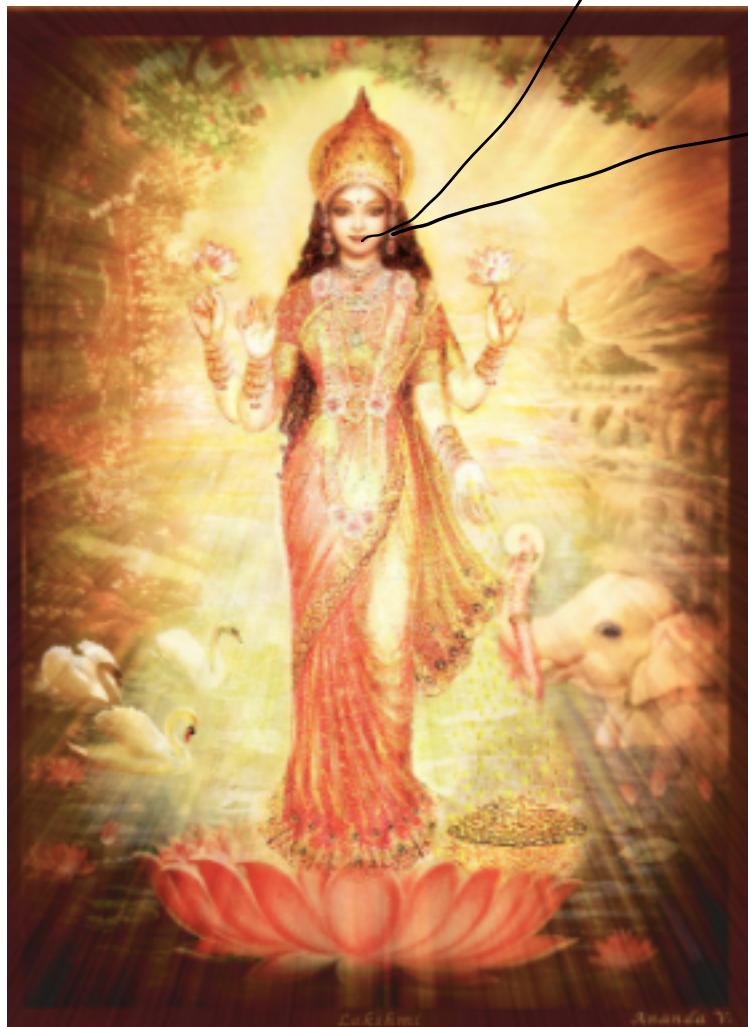


We wish to find the best fitting line through our data (regression line)

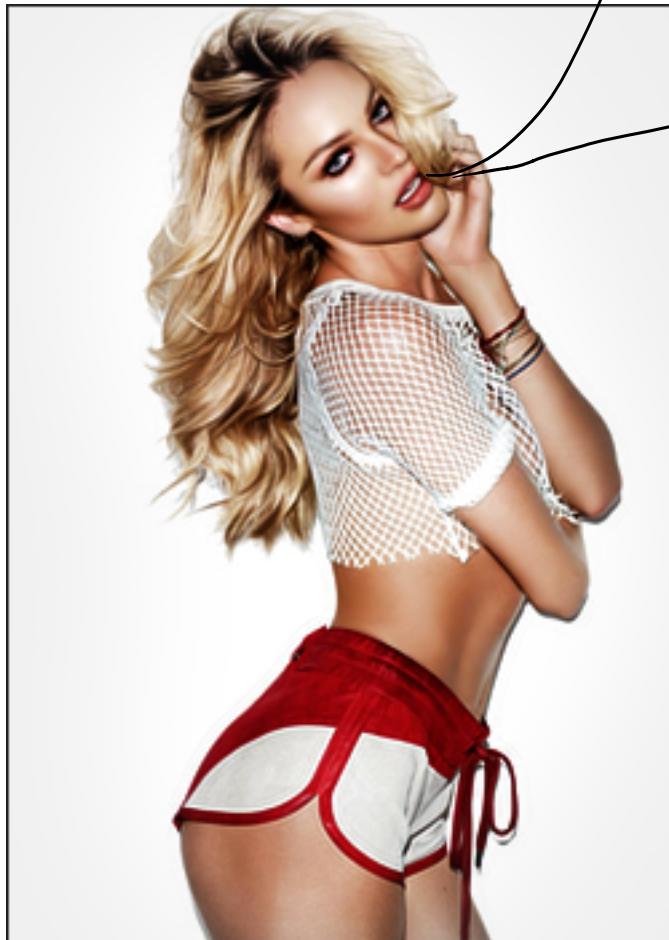
$\hat{e}_i = y - \hat{y}_i$ are the residuals

Since e is homoscedastic, the residuals \hat{e} should also be error homoscedastic

Image of Tyche the goddess of fortune.
Only Tyche knows β_0, β_1 .

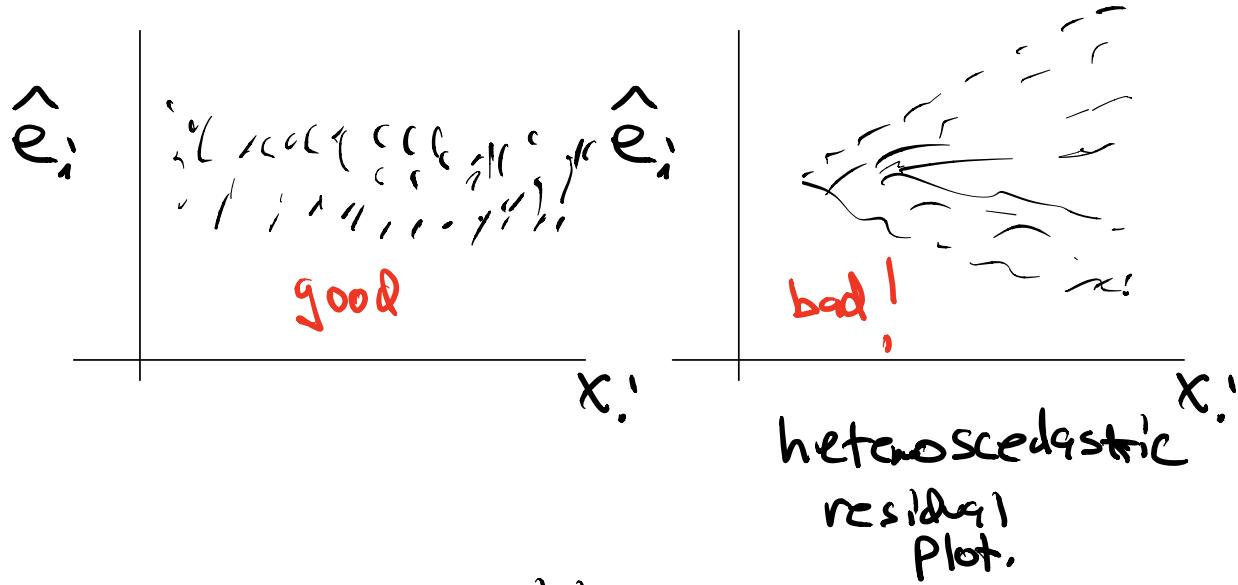


An image of a modern day Tyche from
Sin City.



I know
Eo, B,
!

You should always graph \hat{e}_i , called a residual plot. The residual spread should be constant



We want to minimize the sqrt of sum of squares of the residuals

$$\sqrt{\sum_{i=1}^n \hat{e}_i^2}$$

\uparrow

$\|\hat{e}\|$

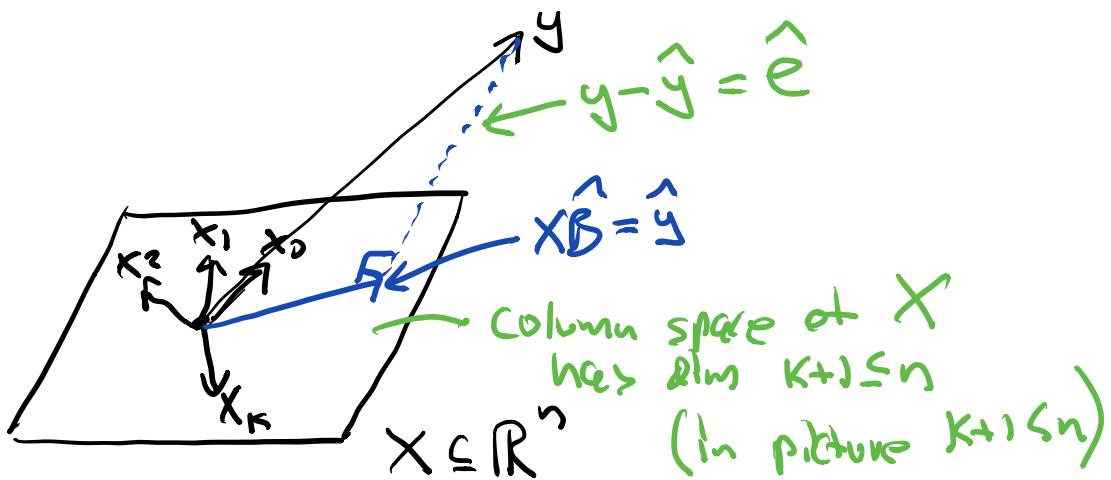
Euclidean norm

of $\hat{e} = \begin{bmatrix} \hat{e}_1 \\ \vdots \\ \hat{e}_n \end{bmatrix}$

③ Geometric derivation of the regression coeffs.

$$y \in \mathbb{R}^n \quad (n \geq k+1)$$

The columns of X span a $(k+1)$ -dimensional subspace since X is full ranks.



Find \hat{B} such that $\|y - x\hat{B}\| = \|y - \hat{y}\| = \|\hat{e}\|$ is as small as possible.

Take the orthogonal projection of y on the column space of X .

$$\underset{(K+1) \times n}{X'} \cdot \underset{n \times 1}{(y - X\hat{B}^{\top})} = \underset{(K+1) \times 1}{0}$$

$$\begin{matrix} x_1 \\ \vdots \\ x_K \end{matrix} \left(\begin{array}{c} \parallel \\ \parallel \\ \parallel \\ \parallel \\ \parallel \end{array} \right) \cdot \underset{n \times 1}{\quad} = \underset{(K+1) \times 1}{\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}}$$

$$\Rightarrow X'y - X'X\hat{B}^{\top} = 0$$

$$\Rightarrow X'y = X'X\hat{B}^{\top}$$

$$\boxed{\hat{B}^{\top} = (X'X)^{-1}X'y}$$

if y lies in column space of X

$$y = X\hat{B}^{\top} \Rightarrow X'y = X'X\hat{B}^{\top}$$

$$\Rightarrow \hat{B}^{\top} = (X'X)^{-1}X'y$$

as
before,

Next time for k=1 case

will show

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} E(y) - \hat{\beta}_1 E(x) \\ \frac{\text{cov}(x,y)}{\text{var}(x)} \end{bmatrix}$$

Last time

Chap 14

Stat 135 Lec 32

Linear Regression

| | | | | |
|-----------|---|---------------------------------------|----------------------|---------------------------------------|
| indep var | cat | cat | quant | quant |
| dep var | quant | cat | quant | cat |
| | t-test (≤ 2 grps) ANOVA (≥ 2 grps) | χ^2 test | linear regression | logistic regression (didn't cover) |
| example | mean wt. of diff. treatment grps. | no. patients of diff. treatment grps. | predict ht. from wt. | predict binary outcome from wt. |
| chapters | 11,12 | 9.5,13 | 14 | |

consider matrices:

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X_0 = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots & 1 \end{bmatrix} \quad X_1 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{n1} \end{bmatrix} \quad \dots \quad X_k = \begin{bmatrix} x_{1k} \\ \vdots \\ x_{nk} \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix} \quad e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$$

$X = \begin{bmatrix} X_0, \dots, X_k \end{bmatrix}_{n \times (k+1)}$

where $n \geq k+1$

and X full rank

design matrix

std statistical model: $y = X\beta + e$

We require e to be homoscedastic (i.e. independent of x).

Properties of linear regression related to the variance of $\hat{\beta}$ relies on this.

Only Tyche knows β and e .

we see the data and can make a

regression "line" $\hat{y} = \hat{\beta}X$ that minimizes the Euclidean norm of the residual vector $\hat{e} = \hat{y} - y$

we found $\hat{\beta} = (X'X)^{-1}X'y$

Today

① Sec 14.3 simple regression

$$\hat{\beta}_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$\hat{\beta}_0 = E(y) - \hat{\beta}_1 E(x)$$

② in-class R exercise

③ Sec 14.4 Statistical properties of $\hat{\beta}$

① Sec 14.3

$$\hat{B} = (\tilde{x}' \tilde{x})^{-1} \tilde{x}' y.$$

In the case to case of simple linear regression
we will show

$$\hat{B}_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$$\hat{B}_0 = \bar{y} - \hat{B}_1 \cdot \bar{x}.$$

$$x = \text{Unif}(x_1, \dots, x_n)$$

$$y = \text{Unif}(y_1, \dots, y_n)$$

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

$$E(y) = \bar{y}$$

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

EXAMPLE B Returning to Example A on fitting a straight line, we have

$$\hat{B}_{2 \times 1} = (\mathbf{x}' \mathbf{x})^{-1} \mathbf{x}' \mathbf{y}$$

where

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 1 & \cdots & 1 \\ x_1 & \cdots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

$$= \begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

Confirm that

$$\hat{B}_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

14.4 Statistical Properties of Least Squares Estimates 567

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{Y} = \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

Thus,

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix}$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix} \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix}$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \left[\begin{array}{c} \left(\sum_{i=1}^n y_i \right) \left(\sum_{i=1}^n x_i^2 \right) - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n x_i y_i \right) \\ n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \end{array} \right] n^2 \text{Cov}(xy)$$

n² Cov(xy) →

which agrees with the earlier calculation. ■

so $\hat{B}_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}$ ✓

We show that $\hat{B}_0 = E(y) - \hat{B}_1 E(x)$.

$$\begin{aligned} \sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i &= \sum y_i \sum x_i^2 - \sum y_i (\sum x_i)^2 - \sum x_i \sum x_i y_i + \sum y_i (\sum x_i)^2 \\ &= \sum y_i \left[\sum x_i^2 - (\sum x_i)^2 \right] - \left[\sum x_i y_i - \sum x_i \sum y_i \right] \sum x_i \\ &= n^2 E(y) \text{var}(x) - n^2 \text{cov}(x, y) E(x) \end{aligned}$$

Hence

$$\begin{aligned} \hat{B}_0 &= \frac{\sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} = \frac{n^2 E(y) \text{var}(x) - n^2 \text{cov}(x, y) E(x)}{n^2 \text{var}(x)} \\ &= E(y) - \frac{\text{cov}(x, y)}{\text{var}(x)} E(x) \quad \text{so } \hat{B}_0 = E(y) - \hat{B}_1 E(x) \quad \checkmark \end{aligned}$$

Next
Find \hat{B} using Calculus approach.

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

To find β_0 and β_1 , we calculate

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum \hat{e}_i = 0$$

$$\sum x_i \hat{e}_i = 0$$

Called
normal
equations

want to
find β_0, β_1 , that
minimizes this.

Setting these partial derivatives equal to zero, we have that the minimizers $\hat{\beta}_0$ and $\hat{\beta}_1$ satisfy

$$\sum_{i=1}^n y_i = n \hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i$$

$$\sum_{i=1}^n x_i y_i = \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2$$

Solving for $\hat{\beta}_0$ and $\hat{\beta}_1$, we obtain

$$\hat{\beta}_0 = \frac{\left(\sum_{i=1}^n x_i^2 \right) \left(\sum_{i=1}^n y_i \right) - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n x_i y_i \right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}$$

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}$$

Problem 10 at the end of the chapter asks you to derive the following useful equivalent expressions:

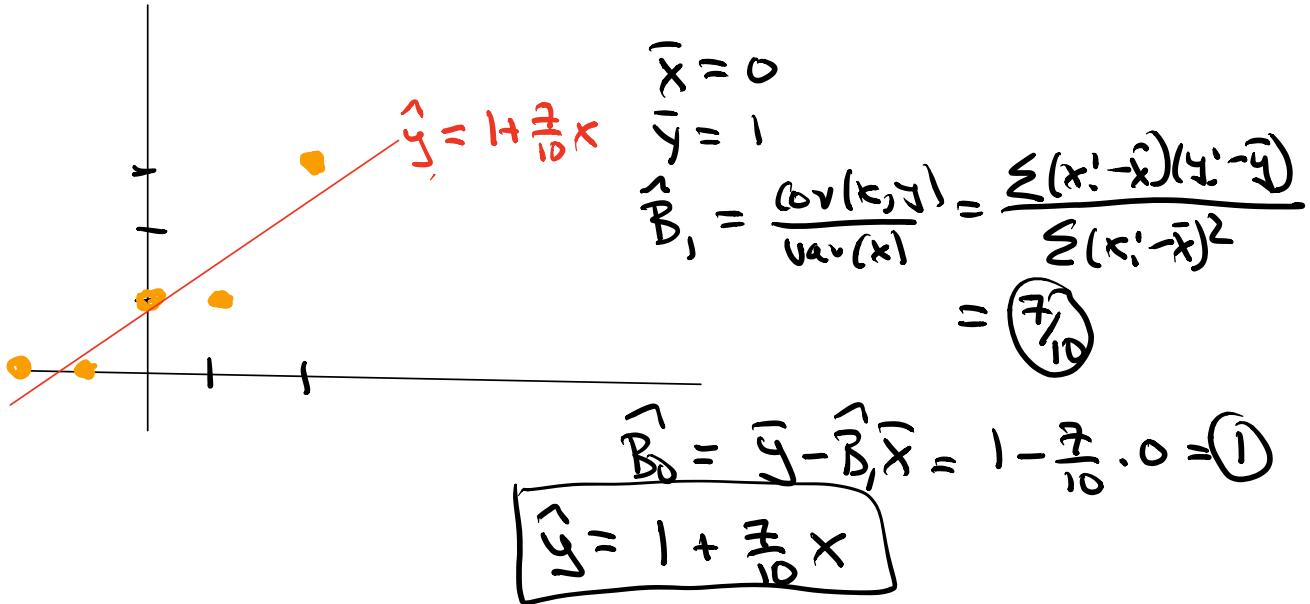
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$E(y) - \hat{\beta}_1 E(x)$

$\frac{Cov(x,y)}{Var(x)}$

Ex Use the method of least squares to fit $n=5$ data pts $(-2, 0), (-1, 0), (0, 1), (1, 1), (2, 3)$



On Hw 9 you will show the regression line can be written

$$\frac{(\hat{y} - \bar{y})}{s_y} = r \frac{(x - \bar{x})}{s_x}$$

Correlation coefficient of $x, y.$



Stat135 lecture 18

Simple Linear Regression

Start Over

In-class exercise
lec 33 on
br courses.

Simple Linear Regression

Here is a basic example of how to draw a linear regression line in R

```
x <- -2:2
y <- c(0,0,1,1,3)
df <- data.frame(x,y)
df
```

← same data as in above example.

| x | y |
|----|---|
| -2 | 0 |
| -1 | 0 |
| 0 | 1 |
| 1 | 1 |
| 2 | 3 |

5 rows

```
reg <- lm(formula=y~x)
summary(reg)
```

← y~x

```
## 
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##   1      2      3      4      5 
## 4.000e-01 -3.000e-01 -2.776e-16 -7.000e-01  6.000e-01 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 1.0000    0.2708   3.693  0.0345 *  
## x           0.7000    0.1915   3.656  0.0354 *  
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
##
## Residual standard error: 0.6055 on 3 degrees of freedom
## Multiple R-squared:  0.8167, Adjusted R-squared:  0.7556 
## F-statistic: 13.36 on 1 and 3 DF,  p-value: 0.03535
```

```
reg$coefficients
```

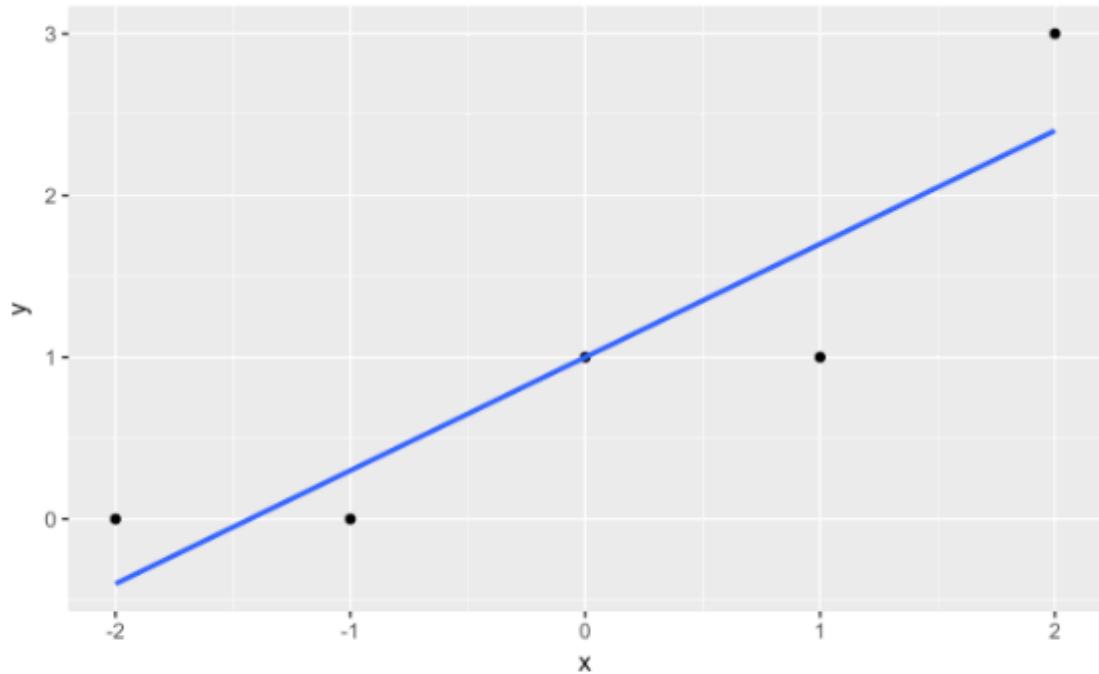
```
## (Intercept)          x
##           1.0           0.7
```

Lets find the expected y value for an x value of 3:

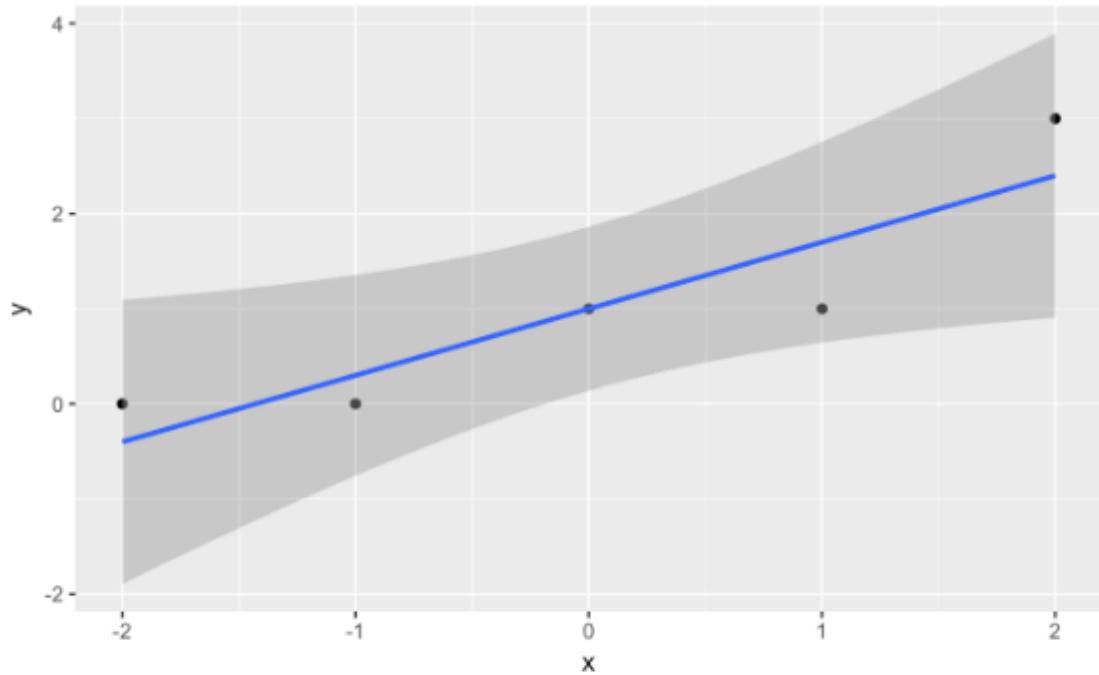
$$\text{reg\$coefficients[1]} + \text{3} * \text{reg\$coefficients[2]} = 1 + 3(0.7)$$

```
## (Intercept)
##           3.1
```

```
df %>% ggplot(aes(x=x,y=y)) + geom_point() + geom_smooth(method=lm,se=FALSE)
```



```
df %>% ggplot(aes(x=x,y=y)) + geom_point() + geom_smooth(method=lm)
```



example of drawing a regression line

There are two variables age and height. ages: 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29
heights: 76.1, 77, 78.1, 78.2, 78.8, 79.7, 79.9, 81.1, 81.2, 81.8, 82.8, 83.5

Find the equation of the regression line predicting height for different ages (i.e. $x=\text{ages}$, $y=\text{heights}$).
Predict the height of a 1 year old. Draw the regression line and the points (no error bars).

Code Start Over Solution
1
2
3

Run Code

Sec 14.4 Statistical properties of least square estimation

A vector-valued RV is called a random vector.

$$\Leftrightarrow \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

We can generalize covariance of a random variable to covariance of a random vector.

Let $\hat{\mathbf{B}} = \begin{bmatrix} \hat{B}_0 \\ \vdots \\ \hat{B}_k \end{bmatrix}$ be random vector

$$E(\hat{\mathbf{B}}) = \begin{bmatrix} E(\hat{B}_0) \\ \vdots \\ E(\hat{B}_k) \end{bmatrix}$$

$$\text{Var}(\hat{\mathbf{B}}) = \sum_{\substack{\text{in R.R.} \\ \hat{\mathbf{B}}\hat{\mathbf{B}}}} \left[\text{Cov}(\hat{B}_i, \hat{B}_j) \right] \quad \begin{array}{l} i=0, \dots, k \\ j=0, \dots, k \\ (k+1) \times (k+1) \end{array}$$

Fact Thm A, B p 533, 574

$$E(\hat{\mathbf{B}}) = \mathbf{B}$$

$$\text{Var}(\hat{\mathbf{B}}) = \sigma^2 (\mathbf{X}' \mathbf{X})^{-1}$$

P574 Rice

E X A M P L E A We return to the case of fitting a straight line. From the computation of $(\mathbf{X}^T \mathbf{X})^{-1}$ in Example B in Section 14.3, we have

$$\Sigma_{\hat{\beta}\hat{\beta}} = \frac{\sigma^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

Therefore,

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{n \text{Var}(x)} \right]$$

algebra see below*

$$\text{Var}(\hat{\beta}_1) = \frac{n \sigma^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \frac{\sigma^2}{n \text{Var}(x)}$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 \sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \frac{-\sigma^2 \bar{x}}{n \text{Var}(x)} ■$$

* algebra:

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

$$\Rightarrow \sum x_i^2 = n\text{var}(x) + n\bar{x}^2$$

$$\Rightarrow \text{var}(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n^2 \text{var}(x)} = \sigma^2 \left(\frac{n\text{var}(x) + n\bar{x}^2}{n^2 \text{var}(x)} \right)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{n \text{var}(x)} \right)$$

Knowing the variances of $\hat{\beta}_0$ and $\hat{\beta}_1$ will allow us to do hypothesis testing on whether x and y are related and make confidence intervals.

We don't know σ^2 so next time we will find an estimator of σ^2 .