

## Chapter 1

**Creating vectors:** c(1,2,3) = 1 2 3, c("a", "b") = "a", "b"  
 $\text{rep}(2,5) = 2 \ 2 \ 2 \ 2 \ 2$   
 $\text{seq}(\text{from} = 2, \text{to} = 10, \text{by} = 3) = 2, 5, 8$   
 $\text{seq}(1, 10, \text{length} = 3) = 1, 0, 5, 5, 10, 0$   
 $\text{class}(a) = \text{"numeric" / "character"}$  (cannot add both tgt)

**Creating matrix:** matrix(1:12, byrow=T, nrow=3)

Matrix multiplication: %\*% Element multiplication: \*  
 $\text{Dim}()$ , +, -, transpose(), inverse(), solve(), determinant: det  
 $\text{diag}(c(1:12))$  diagonal 1,2 others 0  
 $[-r,-c]$ : remove row and col  
 $r/cbind(1:3, c(2,3,5))$ : add row/col

**Creating dataframes:**

```
Customer = c('Ian', 'John', 'Keegan')
Married = c(TRUE, FALSE, TRUE)
Age = c(35, 26, 48)
Opinion = c('Excellent', 'Not bad', 'Good')
```

```
survey = data.frame(Customer, Married, Age, Opinion)
```

```
Survey Subject=rep(1001:1027, rep(3,27))
Customer Married Age Opinion
Treatment=(c("placebo", "Low Dose", "High Dose")
Period=(1:3,3,1,2,2,1)
print(cbind(Subject, Treatment, Period))
```

survey[2,]  $\Rightarrow$  2<sup>nd</sup> row, survey[,3];  $\Rightarrow$  3<sup>rd</sup> column  
subset(survey, Age>=30)

survey[survey\$Age>=30 & survey\$Sex>=2,]

data[order(data\$Var, decreasing=F), ]

```
painscore<-c(0, 10, 0.8)
painvector<-c()
for(i in 1:length(painscore)){
  painvector[i]=cainvein(painscore[i])
}
print(pain_level)
```

```
ifelse(a % 2 == 0, "even", "odd")
```

C= function(n) { value = 0 }

i = 2  
while(i <= n){ value = value + i^3; i=i+1}

return(value)}

```
best3_mean <- function(score) ly(score/order(score, decreasing = TRUE))
average <- best3_mean(score)
return(means(1:3)))
```

best3\_ave <- apply(CA\_result[,2:6],1,best3\_mean)

## Chapter 2

sum(data)/length(data) == mean(data) mean(data, trim=0.1)  
# -10% each end

hist(r) (lg shift mean right, exp left)  
xpt<-seq(-001,100,by=0.1)  
n\_den<-dnorm(xpt,mean(r),sd(r))  
ypt<-n\_den\*length(r)\*10  
lines(xpt,ypt,col="blue")

See if data is normally distributed

1. q=seq(from=0.02,to=0.98,by=0.02)

x=2\*3\*pnorm(q)  $\Rightarrow$  N(2,32)

OR x=q(t,qf=1)  $\Rightarrow$  t distribution v=1

qnorm(x)

qqline(x, col="blue")

2. Shapiro.test(return)  $\Rightarrow$  small values not normal, larger sample more stat significant

boxplot(data)  $\Rightarrow$  Identify outliers

Classic technique: abs(data - mean(data)) > 2\*sd(data)  
Boxplot: (data < quantile(data, 0.25))-1.5\*IQR(data) | (data > quantile(data, 0.75))+1.5\*IQR(data))

median and mean  $\Rightarrow$  boxplot, bigger spread (var)  $\Rightarrow$  hist

## Chapter 3

CLT: for  $n > 30$   $\bar{X} - \mu \sim N(0, 1)$  as  $n \rightarrow \infty$

• When  $\sigma$  known,

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

Pr $\left(-z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_{\alpha/2}\right) = 1 - \alpha$   $\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$

• When  $\sigma$  unknown, use sample sd  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$ , then

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-2, n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-2, n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \sim t_{n-1}$

$\bar{X} - t_{n-1} \cdot \frac{S}{\sqrt{n-1}} < \mu < \bar{X} + t_{n-1} \cdot \frac{S}{\sqrt{n-1}}$

$T$

Using the formula above, we have  
 $S^2 = 98.53$ , and  $T = 25.43$

**Chapter 7**  
**Correlation coefficient:** measure strength and direction of linear relationship (strong/weak/no) between 2 variables

0: no linear relation, strong: 0.8-1, moderate: 0.4-0.8, weak: 0-0.4

$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}} = \frac{E((X - \bar{X})(Y - \bar{Y}))}{\sqrt{E(X - \bar{X})^2}\sqrt{E(Y - \bar{Y})^2}}$  If  $p=1$ , inverse  
 $\rho < 0$ , negative correlation,  $\rho > 0$ , positive correlation

$x = c(10.4, 10.8, 11.1, 10.2, 10.3, 10.5, 10.8, 11.0, 10.6, 11.4)$   
 $y = c(7.4, 7.6, 7.9, 7.2, 7.4, 7.1, 7.4, 7.2, 7.8, 7.3, 7.5, 7.3)$

`> plot(x, y, xlab = "Wing Length", ylab = "Tail Length")` scatter plot

**Kruskal-Wallis rank sum test**

`data: cornfield and cornmethod`  
`Kruskal-Wallis chi-squared = 25.629, df = 3, p-value = 1.14e-05`

**Dependent samples**

**1) Sign Test:** form of quantile test with  $p=0.5$   
 Tests whether one RV in a pair ( $X_i, Y_i$ ) tends to be larger than the other variable.

- "+" if  $X_i < Y_i$ , "-" if  $X_i > Y_i$ , "0" if  $X_i = Y_i$  (discard tied pairs)
- $H_0: P(+)=P(-) = 0.5$ ,  $H_1: P(+) \neq P(-)$
- For  $H_0$  # of '+' should be  $\frac{1}{2}$  sample size  $T = \# \text{ of } + -$
- $T \sim B(n, \frac{1}{2}), E(T) = n/2, \text{Var}(T) = n/4 \Rightarrow T \sim N(n/2, n/4)$

Suppose at the end of the period, 8 preferred B, 1 preferred A, and 1 reported no preference.  
 $\Pr(T=8|T=8, H_0) = 0.0195$   
 $\Pr(T=8|T=8, H_1) = 0.0195$   
 $\Pr(T=8|T=8, \text{one-sided probability}) = 0.0195$   
 $\Pr(T=8|T=8, \text{two-sided probability}) = 0.0390$

So,  $n = 9, T = 8$ . The p-value is  $0.0390$

**Hypothesis Testing**

**Case1:**  $H_0: \rho = 0, H_1: \rho \neq 0$

1: T-test:  $T = \frac{R}{S_R} \text{ df} = n-2, P(|T|) \geq .877, 1.156 = \text{p-value}$

2: F-test:  $F = \frac{1-|R|}{1-|R|} \text{ df} = (n-2, n-2)$

`> cor.test(wt, t)`

can only use t and F distribution for testing  $H_0: \rho = 0, <= 0, >= 0$

**Case2:  $H_0: \rho = \text{any value not } 0$**  eg,  $(H_0: \rho = 0.75, H_1: \rho \neq 0.75)$

Fisher transformation  $r = \frac{1}{2} \ln \left( \frac{1+R}{1-R} \right), z_0 = \frac{1}{2} \ln \left( \frac{1+\rho_0}{1-\rho_0} \right)$

$Z = \frac{z - z_0}{\sqrt{1/(n-3)}} \sim \mathcal{N}(0, 1)$

$z < 0.5 \log((1+r)/(1-r))$ ;  $z > 0.5 \log((1+75)/(1-75))$ ;  $z_0 < 0.5 \log((1+75)/(1-75))$

**Confidence Interval for Correlation**

95% CI for  $z$ :  $z \pm qnorm(1 - 0.025) * (1 / \sqrt{n-3})$

95% CI for  $p$ :  $(exp(z^2/2) - 1) / (exp(z^2/2) + 1) \pm qnorm(1 - 0.025) * (1 / \sqrt{n-3})$

Confidence Interval for Correlation Coefficient

For example, in previous example, we have  $r = 0.870, z = 1.33$  and  $\sigma_z = \sqrt{1/(n-3)} = 1/\sqrt{12} = 0.333$ .

A 95% CI for  $\rho_0$  will be  $\rho_0 \pm 0.25 \times \sigma_z = 0.677 \pm 0.193$

Hence,  $\rho_0 \pm 0.25 \times \sigma_z = 0.677 \pm 0.193$

Lower bound:  $\rho_0 = \exp(2 \times 0.677) - 1 = 0.5896$ , and  $\rho_0 = \exp(2 \times 0.677) + 1 = 0.9504$

Upper bound:  $\rho_0 = \exp(2 \times 1.33) - 1 = 0.9628$ , and  $\rho_0 = \exp(2 \times 1.33) + 1 = 0.9998$

**Simple Linear Regression (X-ind., Y-dependent)**

**Eg. Higher midterm  $\rightarrow$  higher end of term scores**  
 Simple linear regression model based on plots  
 The plots suggest that the distribution of the residuals has slightly longer tails at both ends. This may be due to the two possible outliers at both ends. Other than that, the distribution is fairly symmetric, hence a simple linear regression model is considered reasonable

Intercept slope error

Model:  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Transform non-linear to linear:  $y = \beta_0 e^{\beta_1 x} \Rightarrow \ln y = \ln(\beta_0) + \beta_1 x$

**Fitted regression line:**  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

**Estimated regression line and residuals**  $e_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 \cdot X_i)$

**00 (intercept parameter) hat**

$\hat{\beta}_0 = \sum_{i=1}^n \left( \frac{1}{n} - \frac{\bar{X}_i - \bar{X}}{S_{XX}} \right) Y_i, E[\hat{\beta}_0] = \beta_0$

$Var[\hat{\beta}_0] = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{S_{XX}} \right)$

$T = \frac{\hat{\beta}_0 - \beta_0}{S \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}}}}, \text{df} = n-2$

Previously, we had  $\hat{\beta}_0 = 3.8296, S_{xx} = 4152.182, \bar{x} = 33.45455$  and  $s^2 = 10.4299$

A 95% CI for  $\beta_0$  can be computed as:

$\hat{\beta}_0 \pm t_{n-2, 0.05} \times S \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}}} = 3.8296 \pm 2.0395 \times \sqrt{10.4299 / 33 + 4152.182}$

$\hat{\beta}_0 \pm 2.0395 \times \sqrt{10.4299 / 33 + 4152.182}$

If we want to check whether  $\beta_0$  is zero, then we may test,  $H_0: \beta_0 = 0$  against  $H_1: \beta_0 \neq 0$

Under  $H_0$ , the test statistics

$T = \frac{\hat{\beta}_0 - \beta_0}{S \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}}}} \sim t_{n-2}$

The observed t-statistic is:

$t = \frac{3.8296 - 0}{\sqrt{10.4299} \times \sqrt{\frac{1}{33} + \frac{33.45455^2}{4152.182}}} = 2.1655$

The p-value:  
 $2 \times \Pr(T \geq t) = 2 \times \Pr(T \geq 2.1655) = 0.0382$

**Inferences:** can be

Sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$

Mean square error:  $S^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{SSE}{n-2}$

**Standard error**  $S_R = \sqrt{\frac{1 - R^2}{n-2}}$

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Inferences: can be

Sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$

Mean square error:  $S^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{SSE}{n-2}$

Standard error:  $S_R = \sqrt{\frac{1 - R^2}{n-2}}$

When given a dataset, an estimate would be  
 $s^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{SSE}{n-2} = \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{n-2}$

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Inferences: can be

Sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$

Mean square error:  $S^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{SSE}{n-2}$

Standard error:  $S_R = \sqrt{\frac{1 - R^2}{n-2}}$

When given a dataset, an estimate would be  
 $s^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{SSE}{n-2} = \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{n-2}$

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Inferences: can be

Sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$

Mean square error:  $S^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{SSE}{n-2}$

Standard error:  $S_R = \sqrt{\frac{1 - R^2}{n-2}}$

When given a dataset, an estimate would be  
 $s^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{SSE}{n-2} = \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{n-2}$

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Inferences: can be

Sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$

Mean square error:  $S^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{SSE}{n-2}$

Standard error:  $S_R = \sqrt{\frac{1 - R^2}{n-2}}$

When given a dataset, an estimate would be  
 $s^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2 = \frac{SSE}{n-2} = \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{n-2}$

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)`

Call:  
`lm(formula = waste$y ~ waste$x, data = waste)`

Coefficients:

(Intercept)	waste\$x	$\hat{\beta}_0$
0.8296	0.0299	$\hat{\beta}_1$

Residuals:  
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Estimated regression line:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

The predicted model is:  
 $\hat{Y} = 3.8296 + 0.9036 X$

With SSE = 323.3273

**Use lm()**

`lm(waste$y ~ waste$x, data = waste)</`