

Bioinformatics 525: Module 2

Introduction to Statistics

Lab Session #1

1. Introduction to R (Power Point Slides)
2. Data entry, calculations and graphics. Enter the following height and weight data in R.

Height	Weight
60	84
62	95
64	140
66	155
68	119
70	175
72	145
74	197
76	150

`height=c(60,62,64,66,68,70,72,74,76)`

`weight=c(84,95,140,155,119,175,145,197,150)`

- a. Derive $BMI = \text{Weight(Kg)} / (\text{Height(m)} * \text{Height(m)})$ Or $BMI = 703 * \text{Weight(lb)} / (\text{Height(in)} * \text{Height(in)})$

`BMI=703*weight/(height^2)`

- b. Calculate the following:

	Weight N=9	Height N=	BMI N=
Mean	140		
Variance	1303.25		
SD	26.10		
Median	145		
Q1	119		
Q3	155		
IQR	36 (155-119)		
Min	84		
Max	197		

```
length(weight)
```

```
9
```

```
mean(weight)
```

```
140
```

```
var(weight)
```

```
1330.25
```

```
sd(weight)
```

```
36.10
```

```
summary(weight)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
84	119	145	140	155	197

- c. Graph the histograms, boxplot, and Q-Q plot for weight, height, and BMI using multi figure format (`par(mfrow=c(3,3))`). Do weight, height, and BMI look normally distributed?

```
par(mfrow=c(3,1)) #This option is to get multi-figure display, 3 rows 1 column
```

```
hist(weight); boxplot(weight)
```

```
qqnorm(weight);qqline(weight,lty=2)
```

- d. Use `Shapiro.test()` to test if weight, height, and BMI are normally distributed. Which variables are not normally distributed?

```
shapiro.test(BMI)
```

Shapiro-Wilk normality test

```
data: BMI  
W = 0.8287, p-value = 0.04317
```

p-value < .05, BMI is not normally distributed.

- e. Calculate the overweight variable: `owt=1` if BMI > 25 and 0 otherwise. What is the number and the percent of subjects who are overweight (BMI > 25)?

```
owt<-1*(BMI>25)
```

```
table(owt);prop.table(table(owt))
```

3. TRIal Of Preventing HYpertention (TROPHY) Study.

We will use data from TROPHY study to apply some of the methods presented during the class.

- **Brief Introduction**

TROPHY was an investigator-initiated study to examine whether early treatment of prehypertension might prevent or delay the development of subsequent incident hypertension.

- **Objective**

The primary objective of the study was to determine whether, in patients with prehypertension, two years of treatment with candesartan (at a dose of 16 mg daily) reduces the incidence of hypertension at the end of the 2 year treatment and at two years after the discontinuation of active treatment.

- **Data Set**

The data set is in text format, TROPHY.csv.

- a. Read TROPHY.csv data in RStudio using "Import Dataset" on the Environment Window.

IMPORTANT: type `attach(TROPHY)` to have the variables accessible for analysis.

```
TROPHY=read.csv("Folder/Subfolder/TROPHY.csv")
attach(TROPHY)
```

"Folder/Subfolder" is the path where TROPHY.csv is located

Type `dim(TROPHY)` to get the number of rows (observations) and the number of columns (variables) for this data set. Type `head(TROPHY)` to look at variables name. The following subset of variables is part of this data.

Variable Name	Code
Smoke	Smoking status at baseline: 1=yes/2=no
Age	Age in years at baseline
BMI	Body Mass Index at baseline
Insulin	Insulin at baseline
Gluc_fast	Fasting Glucose at Baseline
Ins_gluc	Insulin:Glucose Ratio at baseline
Triglyceride	Triglyceride at baseline
HDL	High Density Lipoprotein Cholesterol
LDL	Low Density Lipoprotein Cholesterol
HDL_LDL	HDL:LDL Ratio at baseline
Cholesterol	Total Cholesterol at baseline

DBP0	Systolic Blood pressure at baseline
SBP0	Systolic Blood pressure at baseline
BMI24	Body Mass Index at 24 months follow-up
DBP24	Systolic Blood pressure at 24 months follow-up
SBP24	Systolic Blood pressure at 24 months follow-up
HT	Hypertension status at 24 months follow-up: 1=yes/0=No
Trt	1=Candesartan/2=Placebo

- b. What is the mean, sd, Median, Q1, Q3, IQR, Min and Max for baseline blood pressure (DBP0) for each treatment group?

```
mean(DBP0[Trt==1]);mean(DBP0[Trt==2])
```

```
sd(DBP0[Trt==1]); sd(DBP0[Trt==2])
```

```
summary(DBP0[Trt==1]); summary(DBP0[Trt==2])
```

- c. Look at the histograms, boxplot and Q-Q plot to see if HDL is normally distributed

```
par(mfrow=c(1,3))
```

```
hist(HDL)
```

```
boxplot(HDL)
```

```
qqnorm(HDL);qqline(HDL)
```

- d. Use Shapiro Wilks test to show that HDL is not normally distributed.

```
shapiro.test(HDL)
```

Shapiro-Wilk normality test

data: HDL

W = 0.9305, p-value = 1.394e-09

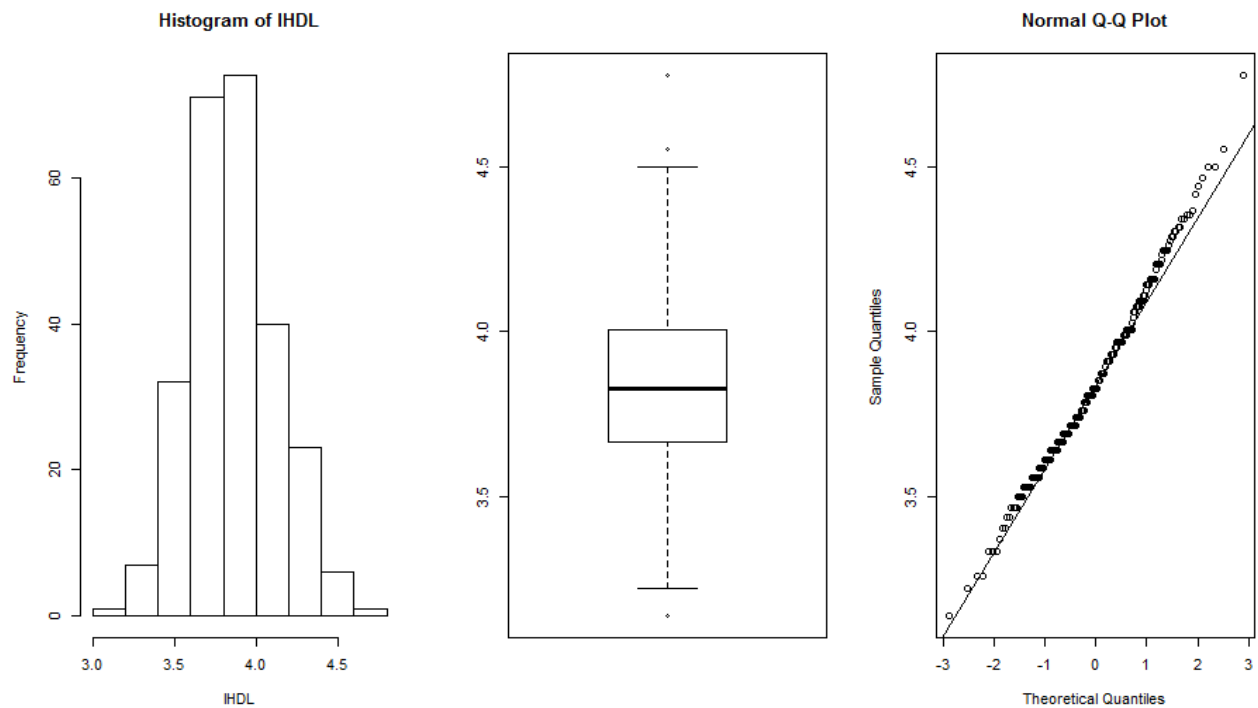
p-value < 0.5, HDL is not normally distributed

- e. Does a log transformation, $IHDL \leftarrow -\log(HDL)$, result in normality for IHDL?

```
par(mfrow=c(1,3))
```

```
hist(IHDL);boxplot(IHDL)
```

```
qqnorm(IHDL);qqline(IHDL)
```



The graphs show no evidence that $\log(\text{IHDl})$ is not normal. The histogram follows a bell curve, the boxplot is symmetric with no severe outliers, and the qqplot is very close to linear.

`shapiro.test(IHDl)`

Shapiro-Wilk normality test

data: IHDl
 $W = 0.9926$, $p\text{-value} = 0.2322$

$p\text{-value} > 0.05$ IHDl is normally distributed.

4. Simulations.

Use simulations in R to illustrate that a linear combination of two normally distributed random variables is normally distributed.

That is, if $z_1 \sim N(m_1, s_1^2)$ and $z_2 \sim N(m_2, s_2^2)$ then $z = a*z_1 + b*z_2 \sim N(a*m_1 + b*m_2, a^2*s_1^2 + b^2*s_2^2)$.

- a. Simulate 1000 data points for $z_1 \sim N(m=1, sd=1)$ and $z_2 \sim N(m=2, sd=2)$

```
z1<-rnorm(1000,m=1,sd=1)
```

```
z2<-rnorm(1000,m=2,s=2)
```

- b. Calculate $z = 3*z_1 + 2*z_2$

```
z<-3*z1+2*z2
```

- c. Use both graphical display tools and the Shapiro Wilks test to test whether z is normal.

```
par(mfrow=c(1,3))
```

```
hist(z);boxplot(z);qqnorm(z);qqline(z)
```

```
shapiro.test(z)
```

- d. Is the mean of z equal to 7(=3*1+2*2); the variance=25(=9*1+4*4); sd=5?

```
mean(z)
```

```
sd(z)
```

- e. Simulate 10000 data points $y \sim N(7,5)$.

```
y<-rnorm(1000,m=7,sd=5)
```

- f. Use summary() function to compare z and y

```
summary(z)
```

```
summary(y)
```

- g. Use side-by-side boxplot to visually compare the distribution of z and y. Are they the same?

```
par(mfrow=c(1,1)) #Changes from the last parameterization
```

```
boxplot(z,y)
```

5. **Simulations.** Use simulations in R to illustrate central limit theorem

- a. Simulate 100 data points for $k=(5,30,100)$ random variable $x_1, x_2, \dots, x_k \sim \text{Bernoulli}(.3)$. The generated data will be a matrix, with k columns (for x_1, x_2, \dots, x_k) and 100 rows.

```
k=100
```

```
x=matrix(rbinom(100*k,1,.3),nrow=100,ncol=k)
```

- b. Calculate $x_{\text{sum}} = x_1 + x_2 + \dots + x_k$

```
xsum=rowSums(x)
```

- c. Look at histogram, Q-Q plot of x_{sum} , does it look normal for different values of k ?

```
par(mfrow=c(1,2))
```

```
hist(xsum);qqnorm(xsum);qqline(xsum)
```

- d. Use Shapiro test to test normality for x_{sum} for different k

```
shapiro.test(xsum)
```