

Week 4: Sampling Distributions, Error and Estimation

D. ODay

2022-07-06

Sampling Error

	Sample	Population
Definitions	Subgroup or portion of the population chosen for evaluation or study	Collection of all items produced or considered
Characteristics	Statistics	Parameters
Size	n	N
Mean	\bar{X}	μ
Median	\tilde{X}	M
Standard Deviation	s	σ
Variance	s^2	σ^2
Skewness	g_3	γ_3
Kurtosis	g_4	γ_4
Proportion	p	π
Rate	\bar{c}	λ

To create random number in R from a specific distribution we use `rnorm()`, `rexp()`, `rpois()`, `rbinom()`.

```
require(lolcat)
```

```
## Loading required package: lolcat
```

```
## lolcat 2.0.0
```

```
# Create four random samples of size n=30 with Mu = 100, sigma = 10
```

```
d1<-rnorm(n = 30,mean = 100, sd = 10)
```

```
d2<-rnorm(n = 30,mean = 100, sd = 10)
```

```
d3<-rnorm(n = 30,mean = 100, sd = 10)
```

```
d4<-rnorm(n = 30,mean = 100, sd = 10)
```

```
# Create a dataframe of all variables
```

```
normdata<-data.frame(d1,d2,d3,d4)
```

```
# Review summary statistics of all variables
```

```
summary.all.variables(normdata, stat.sd=T)
```

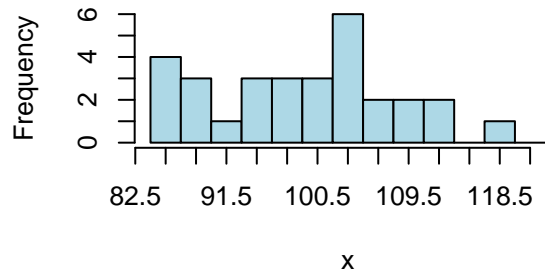
```
##   dv.name  n missing      mean      var      sd g3.skewness  g3test.p
## 1    d1 30      0 99.25340  85.03566  9.221478  0.02639891 0.9476824
## 2    d2 30      0 100.38761  70.28845  8.383821  0.88827060 0.0414139
## 3    d3 30      0 103.50248 142.31778 11.929702  0.46562800 0.2587733
## 4    d4 30      0  98.03228 125.02177 11.181314  0.02495370 0.9505423
##   g4.kurtosis  g4test.p
## 1  -0.6960913 0.35372935
## 2   1.7876385 0.07333562
## 3   1.3310549 0.13476533
## 4  -0.5099145 0.57744477
```

```
# Make output easier to read
nqtr<-function(x,d){noquote(t(round(object(x, d))))}
nqtr(summary.all.variables(normdata, stat.sd=T),3)
```

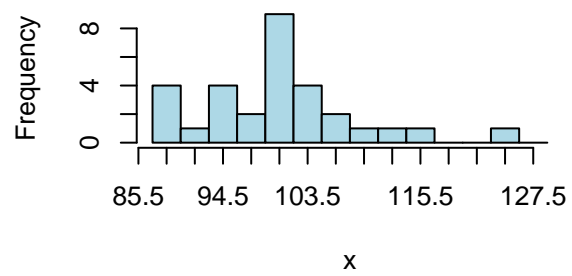
```
##           1      2      3      4
## dv.name   d1     d2     d3     d4
## n         30     30     30     30
## missing   0      0      0      0
## mean      99.253 100.388 103.502 98.032
## var       85.036 70.288 142.318 125.022
## sd        9.221  8.384 11.930 11.181
## g3.skewness 0.026  0.888  0.466  0.025
## g3test.p    0.948  0.041  0.259  0.951
## g4.kurtosis -0.696  1.788  1.331 -0.510
## g4test.p    0.354  0.073  0.135  0.577
```

```
# Create histograms of each variable in one plot
# Set parameters for graphical output, and create a matrix of nrows x ncols
par(mfrow=c(2,2))
hist.grouped(normdata$d1)
hist.grouped(normdata$d2)
hist.grouped(normdata$d3)
hist.grouped(normdata$d4)
```

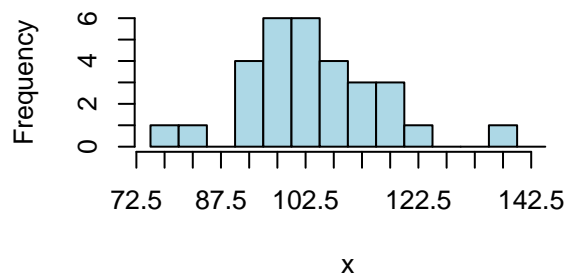
Grouped Histogram



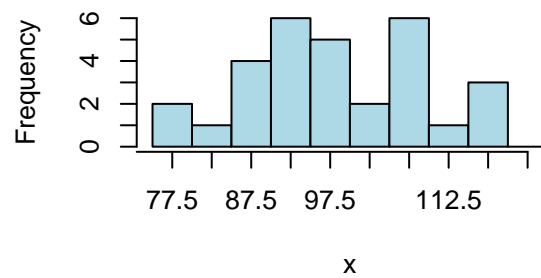
Grouped Histogram



Grouped Histogram

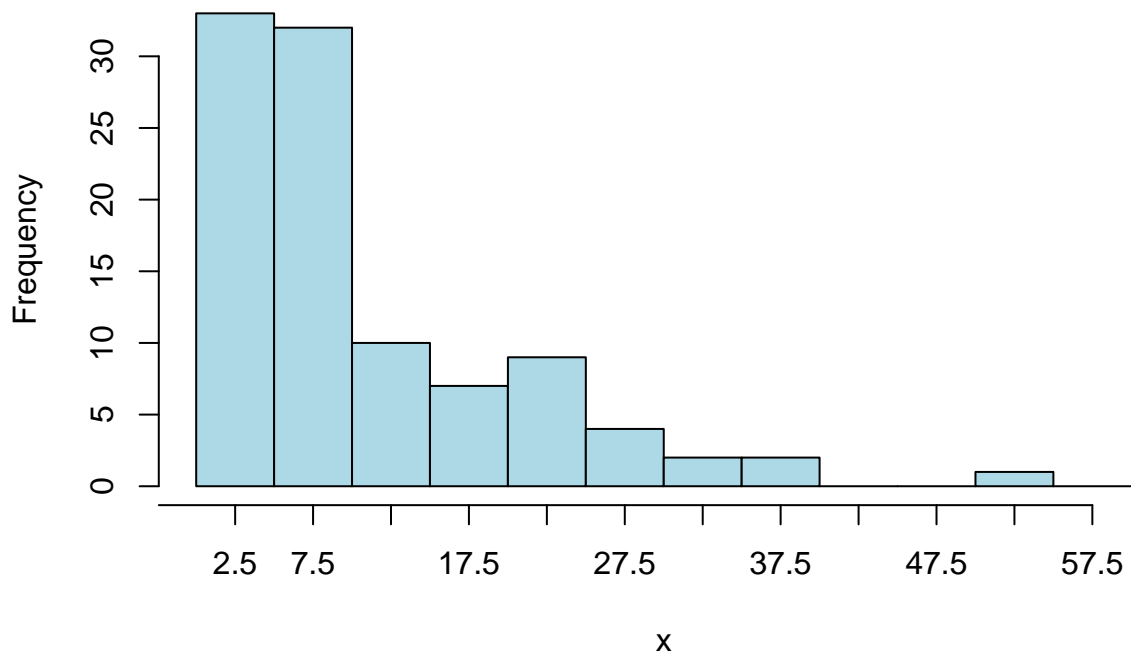


Grouped Histogram



```
# Set parameters back to one graph
par(mfrow=c(1,1))
randoexp<-rexp(n = 100, rate = 1/10)
hist.grouped(randoexp)
```

Grouped Histogram



Sampling Error

- An expected phenomenon since we are not measuring all of the subjects or units for the entire population
- Statistical Methods allow us to account for sampling error, and make appropriate decisions
- In spite of the presence of sampling error, random sampling allow us to use sample statistics as point estimators of population parameters; however, even when unbiased. It will not exactly equal the true value
- An observed difference between a true parameter value and its associated sample descriptive statistic is caused by **sampling error**
- The expected and quantifiable discrepancy between a population parameters and its associated descriptive statistic due to the sample size employed, and in the case of some descriptive statistics, the variability of the population.

-Sampling error is quantifiable using Random Sampling Distributions (RSDs)

- These distributions, like all probability distributions, are based on the principles of classical probability

Random Sampling Distributions and CLT

Random Sampling Distributions

- A RSD is the distribution of a sample statistic calculated from all possible random samples of the given (fixed) size from a given population
 - It is a population distribution and foundational for understanding statistical inference
- To create a RSD
 - Draw all possible random samples of size n from a given research population
 - Calculate descriptive statistics for each of the samples
 - Construct a distribution for each of the sampled descriptive statistics
 - Each of the resultant distributions constitutes the RSD of the statistics

```
# Random Sampling Distribution Simulation
nqtr<-function(x,d){noquote(t(round(object(x, d))))}

# First, set seed so we will get the same results
set.seed(133)

# Create a distribution with mean of 10 and standard deviation of 2
pop<-rnorm(n = 50000, mean = 10, sd = 2)

# Create a histogram of the population distribution
hist.grouped(pop, anchor.value = 0)
hist.grouped(pop, anchor.value=0)
```

RSD of the Sample Averages

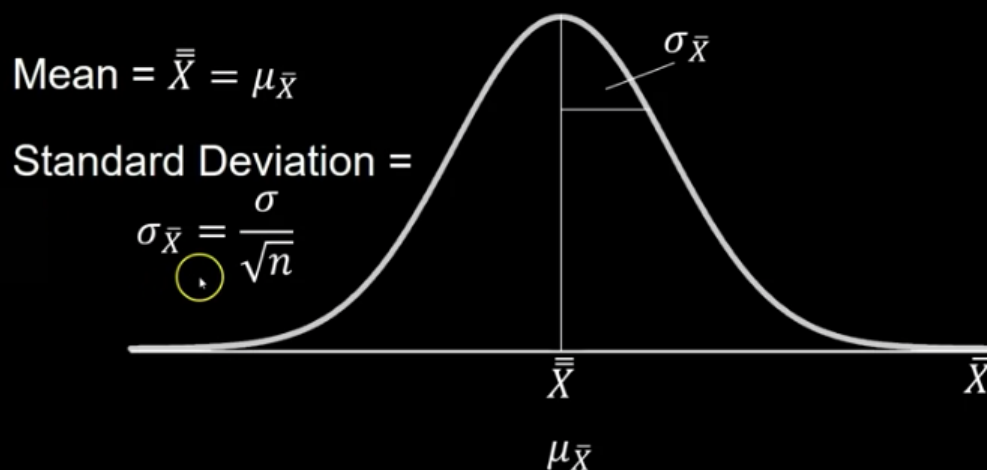
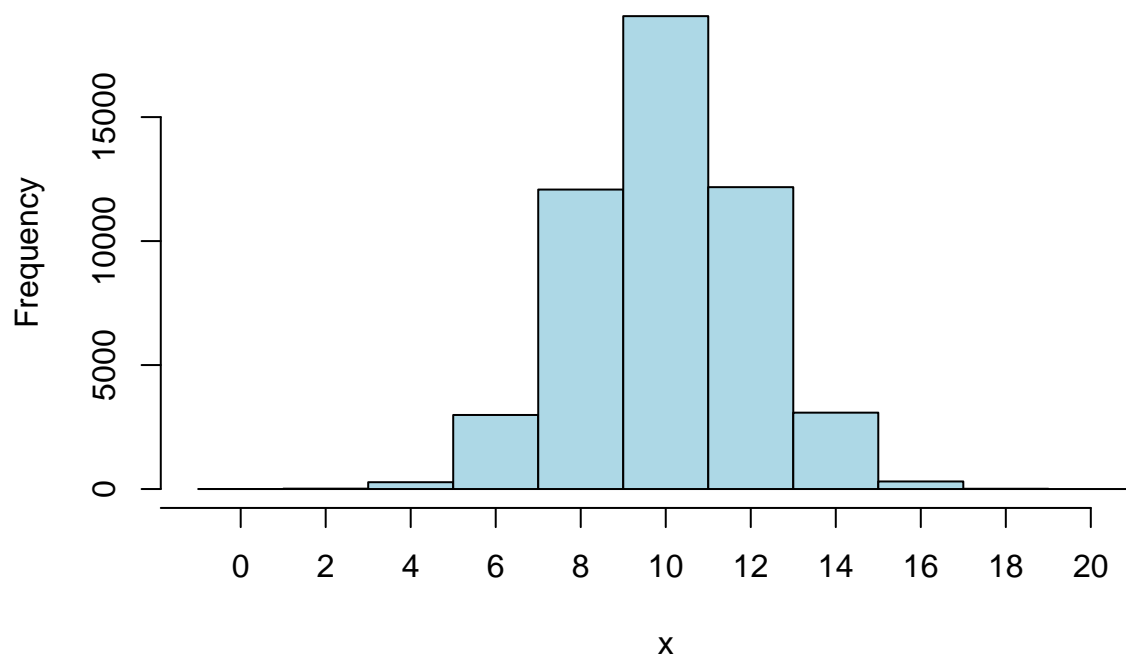


Figure 1: Ex Distribution

Grouped Histogram



```
# Calculate the descriptive statistics of the population distribution
nqtr(summary.continuous(pop, stat.sd=T),5)
```

```
##          1
## dv.name  fx
## n        50000
## missing  0
## mean     10.00851
## var      4.00936
## sd       2.00234
## g3.skewness 0.00476
## g3test.p 0.66361
## g4.kurtosis -0.00231
## g4test.p 0.92471
```

```
# Next, create / simulate a random sampling distribution
# Set the sample size equal to 4
n<-4

# Create the number of repetitions
# (number of times we will take a sample of size 4)
reps<-5000

# Take the random samples from a population with a mean of 10 and
# standard deviation of 2
samples <- replicate(reps, rnorm(n, mean = 10, sd = 2))

# Calculate the averages of each sample of 4
sampleavgs <- colMeans(samples)

# Calculate the descriptive statistics of the RSD
nqtr(summary.continuous(sampleavgs, stat.sd=T),5)
```

```
##          1
## dv.name  fx
## n        5000
## missing  0
## mean     9.97972
## var      1.02914
## sd       1.01447
## g3.skewness 0.01886
## g3test.p 0.58577
## g4.kurtosis -0.01807
## g4test.p 0.81852
```

```
# Create 2 histograms of the population and sample averages
# using the same axis to compare

# Create a 1 x 2 matrix
dev.off() # turns off the default
```

```
## null device
##          1
```

```

layout(matrix(1:2, nrow=2))

# Create both histograms
hist.grouped(pop, xlim=c(0,22), xaxt='n', width.consider = 1)
axis(side = 1, at = seq(0,22,2), labels = seq(0,22,2))

hist.grouped(sampleavgs, xlim=c(0,22), xaxt='n', width.consider = 1)
axis(side = 1, at = seq(0,22,2), labels = seq(0,22,2))

```

Probability with RSDs

Can estimate probability with z-score

$$Z_{\bar{X}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

```

# Using the RSD to solve probability problems
# Example 1
# Define the variables in the problem
mu1<-1.325
sigma1<-0.045
n1<-25
xbar1<-1.433
stderror1<-sigma1/sqrt(n1)

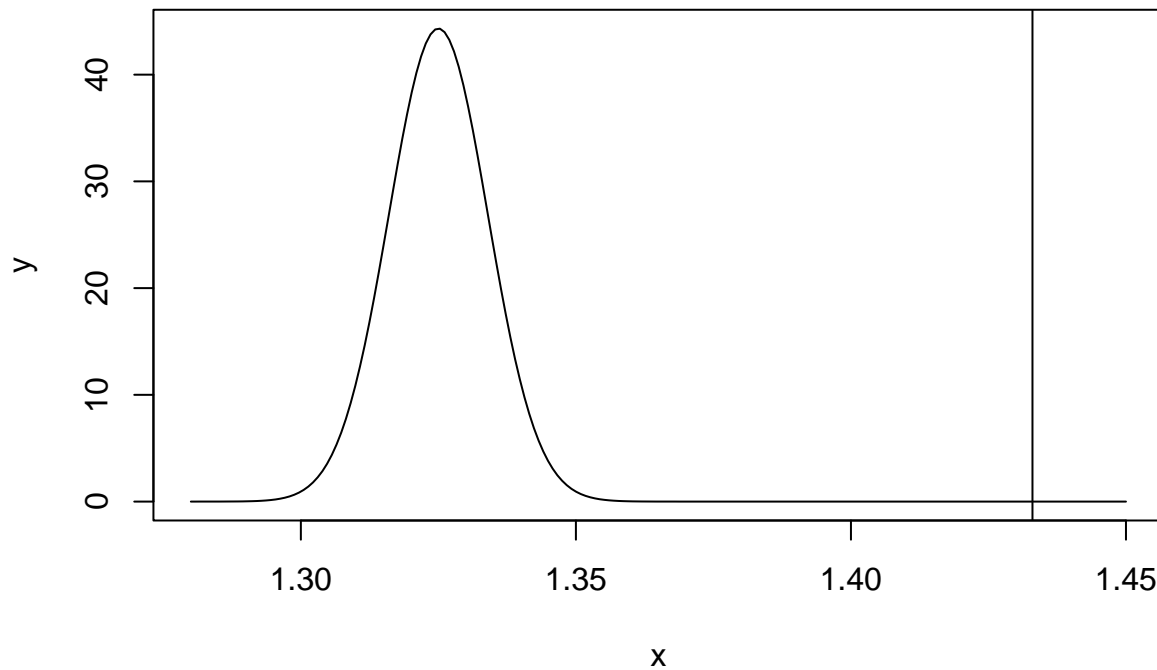
# Calculate the area under the normal curve using the pnorm function
pnorm(q = xbar1, mean = mu1,sd = stderror1, lower.tail = F)

## [1] 1.776482e-33

# Create the normal curve
x=seq(1.28,1.45,length=200)
y=dnorm(x,mean=mu1,sd=stderror1)
plot(x,y,type="l")

# Indicate the location of the xbar of 1.433
abline(v=1.433)

```



```
# Example 2
# Define the variables in the problem
mu2<-50
sigma2<-14.4
n2<-25
xbar2<-55
stderror2<-sigma2/sqrt(n2)

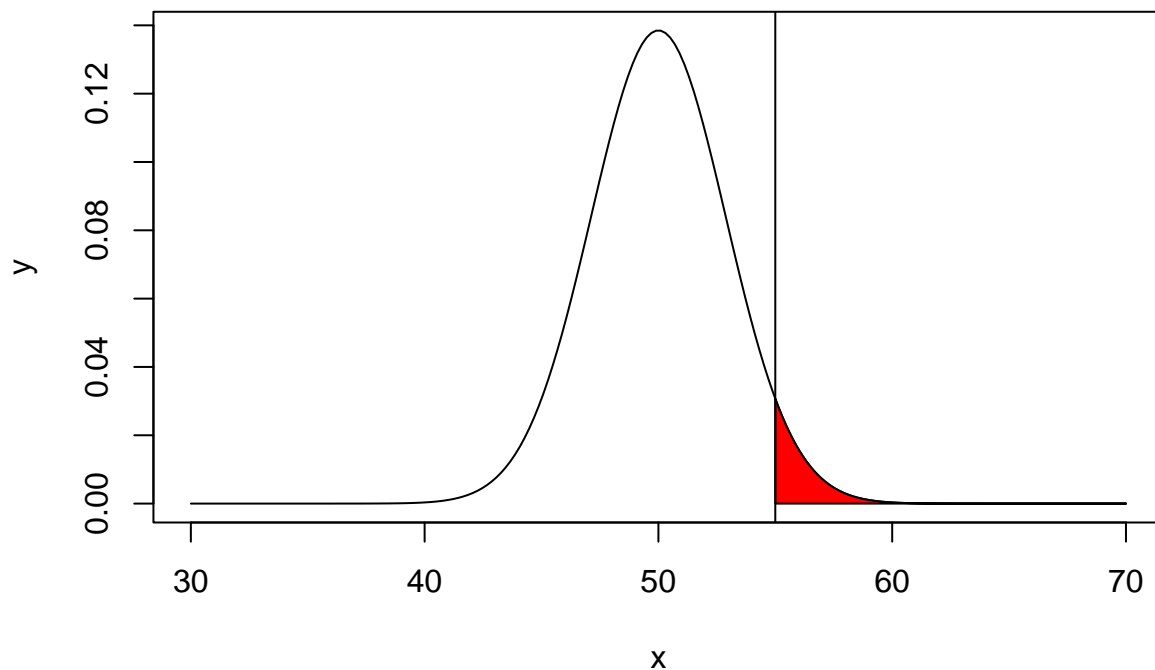
# Calculate the area under the normal curve using the pnorm function
pnorm(q = xbar2, mean = mu2,sd = stderror2, lower.tail = F)
```

```
## [1] 0.0412721
```

```
# Create the normal curve
x=seq(30,70,length=200)
y=dnorm(x,mean=mu2,sd=stderror2)
plot(x,y,type="l")

# Shade the upper tail area
x=seq(55,70,length=100)
y=dnorm(x,mean=mu2,sd=stderror2)
polygon(c(55,x,70),c(0,y,0),col="red")

# Indicate the location of the xbar of 55
abline(v=55)
```

Statistic	RSD	Standard Error
\bar{X}	RSD of the mean	of the mean
\tilde{X}	RSD of the median	of the median
p	RSD of the proportion	of the proportion
R	RSD of the range	of the range

Figure 2: RSD Example

Estimates and Estimators

- **Point Estimate**
 - single number used to estimate an unknown parameter
- **Interval Estimate**
 - Range of values to estimate a parameter
- **Estimator**
 - A sample statistic used to estimate a population parameter. An estimate is a specific observed value of a statistic
- **Criteria for “Good” Estimators**
 - Unbiased

- * the mean of the RSD of the estimator is equal to the parameter it estimates
- Efficiency
 - * The standard error of the statistic RSD. The most efficient estimator is the one with the smallest standard err
- Consistency
 - * Refers to the assumption that as n increases, the value of the statistic approaches the value of its associated population parameter
- Sufficient
 - * Using all possible info the in the sample to estimate the corresponding parameter
 - * example is using sd or var compared to the range since the range only uses two values

Point Estimate		Population Parameter	
Sample Mean	\bar{X}	Population Mean	μ
Sample Variance	s^2	Population Variance	σ^2
Sample Proportion	p	Population Proportion	π
Sample Count	c	Population Count	λ
Sample Skewness	g_3	Population Skewness	γ_3
Sample Kurtosis	g_4	Population Kurtosis	γ_4

Point Estimate	In RStudio
\bar{X}	<code>mean()</code>
s	<code>sd()</code>
p	<code>mean()</code> # average proportion
c	<code>mean()</code> # average count per unit

Confidence Intervals

Confidence Intervals

- The confidence level is the probability associated with an interval estimate
- This refers to the probability that the interval estimate includes the population parameter
- levels are 90, 95 and 99%

-The confidence interval is the range of the estimate. The range of values to find the true population parameter with a given level of confidence

- The interval estimate provides a way to qualify our estimate by indicating the magnitude of the sampling error and hence, the precision of our estimate
- To find this interval, we must look at the set of all possible parameters and assess each of those parameters for their probability of providing us with the sample statistic we observed

Confidence Intervals for the Mean and Variance

```
require(lolcat)
Point_Estimates <- read.delim("~/Documents/GitHub/school_cu/school_cu/methods for quality improvement/D")

ro<-round.object
nqtr<-function(x,d){noquote(t(round.object(x, d)))}
options(scipen=999)
```

- Mean
 - When sd is known
 - When sd is unknown
- Standard Deviation / Variance

Mean (Sigma Known)

$$\mu_{CI} = \bar{X} \pm z_{\alpha/2} * \frac{\sigma}{\sqrt{n}}$$

```
# Confidence Interval for the Mean (Sigma is Known) -----
# Example 1
n<-150
xbar<-20
sd<-5
conf<-0.95

#Use .simple syntax if not using dataframe. Otherwise, don't use the .simple
z.test.onesample.simple(sample.mean = 20
                        ,known.population.variance = 5^2
                        ,sample.size = 150
                        ,conf.level = 0.99)

##
## One-Sample Z Test For Means
##
## data: sample mean, population variance, sample size
## z statistic = 48.99, null hypothesis mean = 0, p-value <
## 0.00000000000000022
## alternative hypothesis: true mean is not equal to 0
## 99 percent confidence interval:
## 18.94842 21.05158
```

```
## sample estimates:
## sample.mean sample.size      se.est      power
## 20.0000000 150.0000000 0.4082483 1.0000000
```

Round the output

```
ro(z.test.onesample.simple(sample.mean = 20
                           ,known.population.variance = 5^2
                           ,sample.size = 150
                           ,conf.level = 0.95),2)
```

```
##
## One-Sample Z Test For Means
##
## data: sample mean, population variance, sample size
## z statistic = 48.99, null hypothesis mean = 0, p-value <
## 0.000000000000000022
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 19.2 20.8
## sample estimates:
## sample.mean sample.size      se.est      power
## 20.00 150.00 0.41 1.00
```

Means (Sigma Unknown)

- Use t distribution
 - lower at the mean and lower at the tails compared to the normal distribution
- t-distribution considers the fact there is error associated with the use of s, the sample standard deviation to estimate the population (s)
- This error increases the variability of the resulting statistic, the t, relative to a standard normal distribution

$$\mu_{CI} = \bar{X} \pm t_{\alpha/2, (n-1)df} * \frac{s}{\sqrt{n}}$$

```
# Confidence Interval for the Mean (Sigma is Unknown) -----
# Example 2
n<-14
xbar<-15000
sd<-500
conf<-0.99

t.test.onesample.simple(sample.mean = xbar
                        ,sample.variance = sd^2
                        ,sample.size = n
                        ,conf.level = conf)
```

```
##
## One-Sample t Test For Means
##
```

```
## data: sample mean, sample size, and estimated variance
## t statistic = 112.25, null hypothesis mean = 0, p-value <
## 0.000000000000000022
## alternative hypothesis: true mean is not equal to 0
## 99 percent confidence interval:
## 14597.47 15402.53
## sample estimates:
## sample.mean      se.est          df var.lowerci      var var.upperci
## 15000.0000      133.6306       13.0000 108989.1895 250000.0000 911632.1111
## sd.lowerci      sd  sd.upperci      power
## 330.1351      500.0000   954.7943      1.0000
```

Standard Deviation

Standard Deviation / Variance

- The following formula may be used to generate a confidence interval for a standard deviation

$$\sqrt{\frac{s^2(n-1)}{\chi_{\alpha/2, (n-1)}^2 df}} < \sigma < \sqrt{\frac{s^2(n-1)}{\chi_{1-\alpha/2, (n-1)}^2 df}}$$

Figure 3: Std Dev Interval

```
# Confidence Interval for the Variance -----
# Example 3
sd<-10
n<-25
conf<-0.95

ci.var<-variance.test.onesample.simple(sample.variance = 10^2
                                         ,sample.size = 25
                                         ,conf.level = 0.95)

# Confidence Interval for the Variance
ci.var$conf.int
```

```
## [1] 60.96929 193.53044
## attr(,"conf.level")
## [1] 0.95
```

```
# Confidence Interval for the Standard Deviation
sqrt(ci.var$conf.int)
```

```
## [1] 7.808284 13.911522
## attr(,"conf.level")
## [1] 0.95
```

```
# Generate Confidence Intervals using a file
# Using the point estimate file, calculate the 95% confidence interval
# estimates for the mean, variance and standard deviation for the Weight data

# Test for normality
summary.continuous(Point_Estimates$Weight)
```

```
## dv.name n missing mean var adtest.AA adtest.p swtest.W swtest.p
## 1 fx 20 0 842.8 19.01053 0.3633424 0.4407723 0.9539818 0.4315902
## g3.skewness g3test.p g4.kurtosis g4test.p
## 1 -0.273123 0.5732011 -0.5046044 0.6871433
```

```
# Calculate confidence intervals
t.test.onesample(Point_Estimates$Weight, conf.level = 0.95)
```

```
##
## One-Sample t Test For Means
##
## data: sample mean, sample size, and estimated variance
## t statistic = 864.46, null hypothesis mean = 0, p-value <
## 0.000000000000000022
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 840.7594 844.8406
## sample estimates:
## sample.mean se.est df var.lowerci var var.upperci
## 842.8000000 0.9749494 19.0000000 10.9946550 19.0105263 40.5545760
## sd.lowerci sd sd.upperci power
## 3.3158189 4.3601062 6.3682475 1.0000000
```

CI for Porportions and Poisson Counts

ex) - $n = 100$ - $p = 0.12$ - Confidence Level Desired = 95%

- $\text{mean} = np = 12$
- $n = 100$
- $a = 0.05$
- $\text{lower_ci} = \text{qbeta}(0.025; 12, 89) = 0.0636$

- $\text{upper_ci} = \text{qbeta}(0.975; 13, 88) = 0.2002$

```
# Confidence Interval for a Proportion -----
# Example 1
ro(proportion.test.onesample.exact.simple(sample.proportion = 0.12
                                           ,sample.size = 100
                                           , conf.level = 0.95),4)
```

```
##
## One-Sample Proportion Test (Exact)
##
## data: sample proportion and sample size
## p = 0.12, null hypothesis proportion = 0.5, p-value <
## 0.000000000000000022
## alternative hypothesis: true proportion is not equal to 0.5
## 95 percent confidence interval:
## 0.0636 0.2002
## sample estimates:
## sample.prop sample.size n.times.p power
## 0.12 100.00 12.00 1.00
```

```
# Using the Point Estimate File, calculate the 90% confidence interval for the
# Proportion
```

```
(prop<-mean(Point_Estimates$Proportion)) # average proportion
```

```
## [1] 0.03525
```

```
(n<-length(Point_Estimates$Proportion)) # sample size
```

```
## [1] 20
```

```
# Proportion (Exact)
ro(proportion.test.onesample.exact.simple(sample.proportion = 0.03525
                                           ,sample.size = 20
                                           ,conf.level = 0.90),4)
```

```
##
## One-Sample Proportion Test (Exact)
##
## data: sample proportion and sample size
## p = 0.0352, null hypothesis proportion = 0.5, p-value <
## 0.000000000000000022
## alternative hypothesis: true proportion is not equal to 0.5
## 90 percent confidence interval:
## 0.0006 0.1949
## sample estimates:
## sample.prop sample.size n.times.p power
## 0.0352 20.0000 0.7050 1.0000
```

Poisson

- $n = 20$
- $\lambda = 25.05$
- desired CI = 95%
- $\lambda * n = 501$
- $\alpha = 0.05$
- $\lambda_{\text{lower}} = \text{qgamma}(0.025; 501)/20 = 22.90$
- $\lambda_{\text{upper}} = \text{qgamma}(0.975; 502)/20 = 27.34$

```
# Confidence Interval for Poisson Counts -----  
# Using the Point Estimate File, calculate the 99% confidence interval for the  
# Poisson Counts
```

```
# Make sure data are Poisson distributed  
poisson.dist.test(Point_Estimates$Count)
```

```
##  
## Poisson Distribution Fit Test Using Variance and Mean  
##  
## data: input data  
## chi.square = 15.926, degrees of freedom = 19, p-value = 0.6756  
## alternative hypothesis: true chi.square is not equal to 19  
## sample estimates:  
##      chi.square sample variance      sample mean  
##      15.92615      20.99737      25.05000
```

```
# Get Total Counts in the Sample  
(counts<-sum(Point_Estimates$Count))
```

```
## [1] 501
```

```
# Get Sample Size  
(n<-length(Point_Estimates$Count))
```

```
## [1] 20
```

```
poisson.test.onesample.simple(sample.count = 501  
                              ,sample.size = 20  
                              ,conf.level = 0.90)
```

```
##  
## Exact Poisson test  
##  
## data: sample.count time base: sample.size  
## number of events = 501, time base = 20, p-value < 0.00000000000000022  
## alternative hypothesis: true event rate is not equal to 1
```



```
## 90 percent confidence interval:  
## 23.23802 26.97066  
## sample estimates:  
## event rate  
## 25.05
```

Homework Quiz

```
set.seed(133)  
d1<-rnorm(n = 10000,mean = 5, sd = 1)  
round(mean(d1), 4)
```

```
## [1] 5.0151
```

```
round(skewness(d1), 4)
```

```
## [1] 0.0378
```

```
round(kurtosis(d1), 4)
```

```
## [1] 0.0411
```

4) Approximately normal