

Function; Variable; Reserved Keywords			done by: Ling 1 Dec 19
	SAS	R	SPSS
Filepath	Filepath = F:\folder\file\	F:/folder/file	
Missing entry	.	NA	.
Comment lines	/* comments */	# comments	-
Importing data	data varname; infile 'filepath' delimiter="," firstobs=2 obs=n; [obs is the final line]	var1 = read.table('filepath', header=TRUE, sep="," col.names, row.names) [attach column names] attach(var2)	"File" -> "Open"
Header names	input col1 col2\$ @@; run;	[col names] names(df1) = c('name1', 'name2') [row names] row.names(df1) = c('name1', 'name2')	*edit directly
Manual entry	datalines; data1 data2 data3 ;	[create vectors] var1 = c(data1, data2, data3); var4 = seq(from=n, to=m, by=x) [create numeric] var2 = numeric(n) [replicates elements] var3 = rep(list, times) [convert vectors to matrix] dim(vector) = c(row, col) [create matrix] var5 = matrix(vector, nrow=n, ncol=m) [create data frames] var6 = data.frame(vector/matrix) [manual entry] var7 = scan()	"Variable View" -> Define variables "Data View" -> Enter data
Fixed format entry	input col1 1-2 col2\$ 4;	var1 = read.fwf('filepath', width=c(1, 2, 3))	"File" -> "Open" -> "Fixed width"
Creating subset	set varname; (filepath if SAS dataset) if condition1 then output;	*reassign variable	"Data" -> "Merge Files"
Filter	if condition1 then output; *note: var must be created first	var[condition,]	[Data View] "Data" -> "Select Cases" -> "If condition is satisfied" -> "If" -> enter conditions
Create new col	col_name = (a+b)/2;	[expand dataframe] var \$col = var [bind a row to matrix] var3 = rbind(var1, var2) [bind a col to matrix] var3 = cbind(var1, var2)	"Edit" -> "Insert Variable"
Keeping subset	keep col1 col2 col3;	var1[rows, cols] *could be logics(conditions) too {note: when filter, apply on rows}	-
Dropping subset	drop col1 col2 col3;		"Edit" -> "Clear"

Function; Variable; Reserved Keywords			done by: Ling 1 Dec 19
	SAS	R	SPSS
Concatenating data set	<code>set var1 var2;</code>	[combining data frames by rows] <code>var3 = rbind(var1, var2)</code>	[Data View] "Data" -> "Merge Files" -> "Add Cases"
If-else	<code>if condition then outcome;</code> <code>else if conditon then output;</code>	<code>if(condition) {statement}</code> <code>else if (condition) {statement}</code> <code>else {statement}</code>	*applicable in filter
Loop	<code>do var = start to end;</code> <code>statements;</code> <code>end;</code> e.g. do batch = 1 to n; do treat = 1 to n; input value @@; output; end; end;	[while loop] <code>while(condition) {statements}</code> [for loop] <code>for (var in 1:n) {statement}</code>	-
Merging data	[Must be sorted by criteria first] <code>merge var1 var2;</code> <code>by criteria;</code>	[merge data by extending cols] <code>var6 =</code> <code>merge(var4, var5, by='criteria', all=TRUE)</code>	[sorted data] "Data View" -> "Data" -> "Merge Files" -> "Add Variables"
Output data	*change the dataname to a permanent address	[save output on screen] <code>sink('filepath') operations sink()</code> [print output] <code>cat(value1, value2, value3, "\n")</code> [save dataframes] <code>write.table(var1, 'filepath', quote=FALSE)</code>	"File" -> "Save"
Sorting	<code>proc sort data=name;</code> <code>by descending coll;</code>	<code>var1[order(criteria),]</code> <code>var2[rev(order(criteria)),]</code> <code>rank(var)</code>	"Data" -> "Sort Cases..."
Writing functions	-	<code>functionname = function(input1, input2) {</code> <code>statement</code> <code>return(output) }</code>	-
Transformation	*refer to create new col	*create new var *e.g. BMI = height/weight**2	[Calculation] "Transform" -> "Compute Variable" [Recode] "Transform" -> "Recode into Different Variable" {note: Click "Output variable are strings" if output values are characters}

Function; Variable; Reserved Keywords					done by: Ling 1 Dec 19																																				
	SAS				R				SPSS																																
Format values	<pre>proc format; *gives lables to values value \$value_name '1' = 'Name1'; value value_name 1-20 = "range" . = "No reply" Other = "Out of range"; data var; label col_name = 'col name' col_name = 'col name'; #label headers format col_name \$ value_name. col_name \$ value_name.;</pre>				<pre>#value labels label = c("label1", "label2") var.group = label[var] #making use of indexing #recode var value = ifelse(test=, yes=, no=) #convert value to index</pre>																																				
Logical Operators	<table><tr><td>equal</td><td>=</td><td>greater than and equal</td><td>>=</td></tr><tr><td>less than</td><td><</td><td>not equal</td><td>^=</td></tr><tr><td>less than and equal</td><td><=</td><td>negation(NOT)</td><td>^</td></tr><tr><td>greater than</td><td>></td><td></td><td></td></tr></table>	equal	=	greater than and equal	>=	less than	<	not equal	^=	less than and equal	<=	negation(NOT)	^	greater than	>				<table><tr><td>equal</td><td>==</td><td>not equal</td><td>!=</td></tr><tr><td>less than</td><td><</td><td>negation(NOT)</td><td>!</td></tr><tr><td>less than and equal</td><td><=</td><td>AND</td><td>&</td></tr><tr><td>greater than</td><td>></td><td>OR</td><td> </td></tr><tr><td>greater than and equal</td><td>>=</td><td></td><td></td></tr></table>	equal	==	not equal	!=	less than	<	negation(NOT)	!	less than and equal	<=	AND	&	greater than	>	OR		greater than and equal	>=				-
equal	=	greater than and equal	>=																																						
less than	<	not equal	^=																																						
less than and equal	<=	negation(NOT)	^																																						
greater than	>																																								
equal	==	not equal	!=																																						
less than	<	negation(NOT)	!																																						
less than and equal	<=	AND	&																																						
greater than	>	OR																																							
greater than and equal	>=																																								
Modify value					<pre>var[row,col] = value</pre>																																				
Applying functions					<pre>list apply, apply function to a list, recoding the list lapply(list, function) table apply tapply(vector, index, function) *the index works like grouping variable apply(array, margin, function) * margin: 1 indicates rows, 2 indicates columns, c(1, 2) indicates rows and columns. Where X has named dimnames, it can be a character vector selecting dimension names.</pre>																																				

Function; Variable; Reserved Keywords		SAS	SPSS
Describing numerical data			
Numerical descriptive measures	<ul style="list-style-type: none"> - Location - Variability - Other measures 	<pre>proc means data=var n mean std min max stderr maxdec=n; title "this is a title"; var col_name; (class col_name; *different groups) OR proc univariate data=var; title "this is a title"; var col_name;</pre>	"Analyze" -> "Descriptive Statistics" -> "Frequencies"
Graphical methods	- Histogram	<pre>proc univariate data=var; var col_name; histogram col_name / midpoints=n to m by 10 normal; class col_name; *seperate by group qqplot col_name;</pre>	"Analyze" -> "Descriptive Statistics" -> "Explore..." OR ~-> "QQ Plots..." "Graphs" -> "Legacy Dialogs" -> "Histogram"/"Box plot" *add <Factor List> for 2 variable analysis
	- Boxplot, Stem and leaf plot, QQ(quantile-quantile) plot	<pre>proc univariate data=var plot; var col_name; OR proc boxplot data=var; plot col_name*col_name;</pre>	
	- Scatter plot for bivariate data	<pre>proc sgplot data=var; title "Title"; by col_name; #different graph scatter x=col_name y=col_name / group=criteria; #same graph</pre>	"Graphs" -> "Chart Builder" *need to drag x and y factors, set colour OR in "Legacy Dialogs"
Robust Location estimator	- Trimmed mean	<pre>proc univariate data=var trimmed=0.2 winsorized=0.2; var col_name;</pre>	-
	- Winsorized mean		-
	- Huber's M-estimators	-	"Analyze:" ->
	- Tukey's bisquare estimator	-	"Descriptive Statistics"
	- Hampel's M-estimator	-	-> "Explore..."
Robust measures of scale parameter	- Interquartile Range (IQR)	<pre>proc univariate data=var robustscale; var col_name;</pre>	
	- Median Absolute Deviation (MAD)		
	- Gini's mean difference		

Function; Variable; Reserved Keywords		SAS	SPSS
Categorical data	- Frequency tables	<code>proc freq data=var; title "Title"; tables col_name -- col_name;</code>	"Analyze" -> "Descriptive Statistics" -> "Frequencies"
	- Contingency tables (two-way frequency table)	<code>tables col_name*col_name;</code>	"Analyze" -> "Descriptive Statistics" -> "Crosstabs..." -> "Stat"
	- Chi-square tests	<code>tables col_name*col_name / chisq; weight col_name; #for manual entry summarised table</code>	"Data" -> "Weight Cases" -> contingency
	- Charts		
Categorical data (paired data)	- McNemar's test	<code>proc freq data=var; tables before*after/agree; weight count; #for frequency counts</code>	"Analyze" -> "Descriptive Statistic" -> "Crosstabs"
	- Bar chart	<code>proc sgplot data=var; vbar col_name; #for two way frequency table panelby col_name; vbar col_name /group=group_name;</code>	"Graphs" -> "Legacy Dialogs.." -> "Bar..."
Numerical data Test for location	Parametric test - One sample t-test (normal population)	<code>proc univariate data=var mu0=value; var col_name;</code>	"Analyze" -> "Compare Means" -> "One Sample t tests..."
	Nonparametric test - Sign test		"Analyze" -> "Nonparametric Tests" -> "Legacy Dialogs" -> "2 Related Samples" *note: must duplicate sample mean in a new col first
	- Sign rank test		
Tests comparing two groups	Two independent sample - Parametric Test (t-test)	<code>proc ttest data=var hu0=value sides=2/L/U; class col_name; #two groups var col_name; #values</code>	"Analyze" -> "Compare Means" -> Independent Sample T tests..."

Function; Variable; Reserved Keywords		SAS	SPSS
	- Nonparametric tests (Wilcoxon rank sum test) (Mann-Whitney U-test)	<code>proc npar1way data=var wilcoxon;</code> <code>class col_name; #two groups</code> <code>var col_name; #values</code> <code>exact wilcoxon;</code>	#recode into numeric variable "Transform" -> "Recode" -> "Into different variable..." #wilcoxon rank sum test "Analyze" -> "Nonparametric Tests" -> "Legacy Dialogs" -> "2 Independent Samples..."
	Related samples - Paired t-test	<code>proc ttest data=var mu0=value sides=2/L/U;</code> <code>paired var1*var2;</code> #or create a difference and run 1 var t-test	"Analyze" -> "Compare Means" -> "paired Sample T test..."
	Related samples: Nonparametric tests - Wilcoxon signed rank test	#create a difference at data import stage <code>proc univariate data=var;</code> <code>var col_name;</code>	"Analyze" -> "Nonparametric Tests" -> "Legacy Dialogs" -> "2 Related Samples..."
Tests comparing three groups or more	Parametric Test - ANOVA table	<code>proc glm data=var;</code> <code>class col_name;</code> <code>model y=x;</code>	"Analyze" -> "Compare Means" -> "One-way ANOVA..." "Options", "Descriptive"
	- LSD	<code>means col_name/ lsd;</code>	"Analyze" -> "Compare Means" -> "One-way ANOVA..." "Post Hoc..", "LSD"
	- Contrast	<code>contrast "X vs. Y and Z" y_col -2 1 1; C: $-2\mu_x + \mu_y + \mu_z = 0$</code> <code>contrast "Y vs. Z" y_col 0 1 -1; C: $\mu_y - \mu_z = 0$</code>	"Analyze" -> "Compare Means" -> "One-way ANOVA..." "Contrast"
	Testing for Equal Variances - Levene's test & Bartlett test	<code>means col_name/ hovtest = levene hovtest = bartlett;</code>	"Analyze" -> "Compare Means" -> "One-way ANOVA..." "Options", "Homogeneity of variance test"

Function; Variable; Reserved Keywords		SAS	SPSS
	Testing for Independence and Normality Assumptions	<pre> *obtain the residuals; proc glm data=var; class col_name; model y=x; output out=var2 p=yhat r=resid; *test for normality; proc univariate data=var; var resid; histogram resid/ normal noplot; qqplot resid; *residual plots; proc sgplot data=var; scatter x=yhat y=resid; </pre>	<p>“Analyze” -> “General Linear Model” -> “Univariate..”</p> <p>“Dependent Variable”, “Fixed Factor”</p>
	Nonparametric 1-way test - Kruskal-Wallis Test	<pre> proc npar1way data=var; class col_name; var col_name; exact wilcoxon; </pre>	<p>“Analyze” -> “Nonparametric Tests” -> “Legacy Dialogs” -> “K Independent Samples...”</p> <p>“Test Variable List”, “Grouping Variable”</p> <p>“Define Range”</p>
Regression Analysis	Pearson Correlation Coefficient	<pre> proc corr data=var nosimple; var col_name1 col_name2 col_name3; </pre>	<p>“Analyze” -> “Correlate” -> “Bivariate...”</p> <p>“Variables window”, “Pearson”</p>
	Simple regression	<pre> proc reg data=var; model y = x; output out=var p=yhat r=resid; quit; </pre>	<p>“Analyze” -> “Regression” -> “Linear”</p>
	Testing assumptions - Linear relationship: scatter plot	<pre> proc sgplot data=var; scatter y=col_name x=col_name; </pre>	<p>“Graphs” -> “Legacy Dialogs...” -> “Scatter/plot”</p> <p>“Simple Scatter” -> “Define”</p>

Function; Variable; Reserved Keywords			SAS	SPSS
	- Normality and independence		<pre>proc sgplot data=var; scatter y=resid x=yhat; proc sgplot data=var; scatter y=resid x=height;</pre>	“Analyze” -> “Regression” -> “Linear” [and save residuals] “Analyze” -> “Descriptive Statistics” -> “QQ plots...”
	- Prediction	Draw fitted lines	<pre>proc sgplot data=var; pbspline y=col_name x=col_name; *pbspline = penalized beta spline method;</pre>	
		Case1: value in dataset	Add statement <code>output out=var p=yhat r=resid</code> and print dataset to search	“Analyze” -> “Regresion” -> “Linear” “Save”, “Unstandardized”
		Case2: value not in	Add data point into dataset with the missing values (using <code>.</code> for missing value)	Add missing data point
	- Transform	Quadratic	<pre>data var2; set var1; x2 = x**2; proc reg data=var; model y=x x2;</pre>	“Transform” -> “Compute Variable...”
		Log	<pre>data var2; set var1; lx = log(x); proc reg data=var; model y=x lx;</pre>	

Function; Variable; Reserved Keywords		R	
Describing numerical data			
Numerical descriptive measures			
- Location	- Mean	<code>mean(var)</code>	<code>summary(var)</code>
	- Median	<code>median(var)</code>	Min/1stQ/Median/Mean/3rdQ/Max
	- Mode	<code>var[order(var)[1:n]]#smallest 5</code>	
- Variability	- Variance	<code>var(var)</code>	
	- Standard deviation	<code>sd(var)</code>	
	- Range	<code>range(var)</code>	
	- Interquartile range	<code>IQR(var)</code>	
	- Coefficient of Variation	<code>cv = function(x) sd(x)/mean(x)*100</code>	
- Other measures	- Min and Max	<code>min(var); max(var)</code>	
	- First & third quartile	<code>quantile(var, level)</code>	
	- Percentiles		
		<code>n = length(x); mn = mean((x-mean(x)^n))</code>	
	- Skewness [measure the symmetric]	<code>skew = function(x) { sk = m3/m2^(3/2)*sqrt(n*(n-1))/(n-2) return(sk)}</code>	
	- Kurtosis [measure the tail]	<code>kurt = function(x) { ku = (n-1)/((n-2)*(n-3))*((n+1)*m4/m2^2-3*(n-1)) return(ku)}</code>	
Graphical methods			
- Histogram	<code>par(mfrow=c(2,1)) #2 graphs in one column in 1 page hist(var, include.lowest=TRUE, freq=TRUE, col='color', main=paste("Title"), sub=paste("sub-title"), xlab="x", ylab="y", axes=TRUE) #normal curve imposed on the histogram xpt = seq(n,m,0.1); ypt=dnorm(seq(n,m,0.1),mean(col_name),sd(col_name)) aypt=ypt*length(col_name)*10 #scaling the line to match area of histogram, 10 is the width of bin lines(xpt,aypt) par(mfrow=c(1,1)) #1 graphs in one column in 1 page</code>		
- QQplot	<code>qqnorm(col_name); qqline(col_name)</code>		
- Boxplot	<code>boxplot(col_name~col_name2) #~means split by group</code>		
- Stem and leaf plot	<code>stem(col_name)</code>		
- Scatter plot for bivariate data	<code>plot(x=col_name,y=col_name, xlim=c(lower,upper), ylim=c(lower,upper), axes=FALSE, pch=0, col=2) par(new=TRUE) #plot new graph and old graph together</code>		
Robust Location estimator			
- Trimmed mean	<code>mean(var, trim=0.2)</code>		
- Winsorized mean			
- Huber's M-estimators			
- Tukey's bisquare estimator			
- Hampel's M-estimator			
Robust measures of scale parameter			
- Interquartile Range (IQR)	<code>IQR(var)/1.34898</code>		
- Median Absolute Deviation (MAD)	<code>median(abs(x-median(x))) mad(x)</code>		
- Gini's mean difference			
Categorical data			
- Frequency tables	<code>table(var)</code>		
- Contingency tables	<code>table(var,var2)</code>		
- Fisher test	<code>fisher.test(table(var), alternative="less/greater/two.sided")</code>		

Function; Variable; Reserved Keywords		R
- Chi-square tests		<code>chisq.test(table(var), correct=F)</code> #applied w/o continuity correction
- Charts		
Paired data	- McNemar's test	<code>mcnemar.test(table1, correct=T)</code>
	- Bar chart	<code>barplot(height= c(table1), names.arg=c("Name1", "Name2"))</code> #two-way frequency table <code>barplot(tablecount, beside=T, col=c(1,8))</code> <code>legend("top-left", c("Name1", "Name2"), fill=c(1,8))</code>
Numerical data		
Tests for location	Parametric test - One sample t-test	<code>t.test(var, mu=value, alternative="two.sided/less/greater")</code>
	Nonparametric test - Sign test	<code>testvalue = var[var!=mu0]</code> <code>binom.test(sum(testvalue>mu0), length(testvalue))</code>
	- Signed Rank test	<code>wilcox.test(testvalue, mu=value, alternative="two.sided")</code>
Two independent samples	- t-test	#first create different groups #test for variances <code>var.test(var1, var2)</code> #t-test <code>t.test(var1, var2, mu=0, var.equal=T/F, alternative="two.sided")</code>
	- Nonparametric (Wilcoxon rank sum test)	#first create different groups <code>wilcox.test(groupa, groupb)</code> or <code>wilcox.test(y~x)</code>
Two related samples	- Paired t-test	<code>t.test(var1, var2, mu=0, paired=T)</code>
	Nonparametric test - Signed test	#create a diff first <code>diff = varA-varB</code> <code>ncount = sum(diff>0)</code> <code>binom.test(ncount, length(diff), 0.5)</code>
	- Signed rank test	<code>wilcox.test(diff)</code>
Three or more independent samples	Parametric test - ANOVA	<code>modell = anov(y~x)</code> <code>summary(modell)</code> <code>tapply(x,y,mean); mean(x); tapply(x,y,mean) - mean(x); boxplot(y~x)</code>
	- LSD	<code>group.1 = y[x=="1"]; group.2 = y[x=="1"];</code> <code>group.3 = y[x=="3"];</code> <code>group.mean = tapply(x,y,mean)</code> <code>treat.group = cbind(group.1, group.2, group.3)</code> #d=(N-k) degree of freedom, where k is the no. of groups <code>mse = sum(modell\$res^2)/d</code> <code>lsd = qt(0.975,d)*sqrt(mse* (1/n1 + 1/n2))</code> <code>diff = mean(treat.group[,i]) - mean(treat.group[,j])</code> <code>if(abs(diff)>lsd) {cat("There is significant difference between groups")}</code>
	- Contrast	<code>contrasts(x) = matrix(c(2,-1,-1,0, -1,1),nrow=3)</code> <code>modelc = lm(y~x)</code> <code>summary(modelc)</code>
	Testing for Equal Variances - Levene's test	<code>group.means = tapply(y,x,mean)</code> <code>gmean = group.means[x]</code> #Absolute deviation from group means is used <code>words.abs.dev = abs(y-gmean)</code> <code>modellev = aov(words.abs.dev~x)</code> <code>summary(modellev)</code>
	- Bartlett's Test	<code>bartlett.test(y,x)</code>

Function; Variable; Reserved Keywords		R
	Test for independence and Normality assumptions	<pre>#Test for normality (Kolmogorov-Smirnov test) resid=model1\$res ks.test(resid,"pnorm",mean(resid),sd(resid)) #Model checking: residual plots fv = model1\$fitted par(mfrow=c(2,1)) plot(fv,resid); abline(h=0,lty=2) qqnorm(resid); qqline(resid,lty=2) par(mfrow=c(1,1))</pre>
	Nonparametric 1-way Test - Kruskal-Wallis Test	<pre>kruskal.test(y,x) or kruskal.test(y~x)</pre>
Simulation Study		
Condition of interest - Compare 3 estimators for location μ through a simulation study - 3 estimators are: sample mean, sample median, sample 10% trimmed mean - Underlying condition: $N(1,1)$ - Sample size: 15 - Simulation size: 1000		
Comparing the estimators	Environment set up	<pre>ns = 1000; n = 15; mu = 1; sd = 1; meax = numeric(ns); medx = numeric(ns); trmx = numeric(ns); stdx = numeric(ns); set.seed(12345)</pre>
	Compute the numeric value of the estimators	<pre>for (i in 1:ns){ x = rnorm(n, mu, sd) meax[i] = mean(x) medx[i] = median(x) trmx[i] = mean(x, trim=0.1) stdx[i] = sd(x) }</pre>
	Compute the mean of estimators	<pre>simumean = apply(cbind(meax,medx,trmx),2,mean)</pre>
	Compute the sd of estimators	<pre>simustd = apply(cbind(meax, medx, trmx), 2, sd)</pre>
	Compute bias	<pre>simubias = simumean - rep(mu,3)</pre>
	Compute MSE	<pre>simumse = simubias^2 + simustd^2</pre>
	Presentation	<pre>sumdata = rbind(c(mu,mu,mu), ns, simumean, simustd, simubias, simumse) col.name = c("mean","median","10% tmean") row.name = c("True value", "No. of simu", "MC Mean", "MC sd", "MC Bias", "MC MSE") dimnames(sumdata) = list(row.name,col.name) round(sumdata,4)</pre>
Checking coverage probability of confidence interval	getting $t_{0.025,n-1}$	<pre>t05 = qt(0.975,n-1)</pre>
	Compute CI	<pre>lower = meax-t05*stdx/sqrt(n) upper = meax+t05*stdx /sqrt(n)</pre>
	Test logic	<pre>coverage = sum(lower <= mu & upper >= mu)</pre>
	Compute coverage	<pre>coverage = coverage/ns</pre>
Checking size of t-test	Calculate t stat	<pre>ttests = (meanx-mu) / (stdx/sqrt(n))</pre>
	Compute rejection	<pre>reject = abs(ttests) > t05</pre>
	Calculate Size	<pre>size = sum(reject)/ns</pre>
Checking power of t-test when $\mu_1 = \mu_0 + \frac{1}{2}\sigma$	Generate data under H_1	<pre>ns = 1000; n = 15; mu0 = 1; sigma = 1; mu1 = mu0 + 0.5*sigma meanx = numeric(ns); stdx = numeric(ns); set.seed(12345) for (i in 1:ns){ x = rnorm(n, mu1, sigma) meanx[i] = mean(x) stdx[i] = sd(x) }</pre>

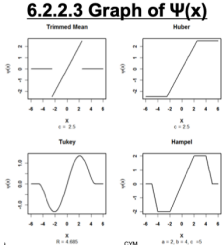
Function; Variable; Reserved Keywords		R
	Calculate t stat	<code>t05 = qt(0.975,n-1)</code> <code>ttests = (meanx-mu0)/(stdx/sqrt(n))</code>
	Compute rejection	<code>#alternate method of calculating rejection</code> <code>reject = ttests >= t05 ttests <= -t05</code>
	Calculate power	<code>power = sum(reject)/ns</code>
Bootstrap method		
Generate bootstrap samples	Compute true value	<code>options(digit=4)</code> <code>theta.hat = corr(var1,var2)</code>
	Generate bootstrap samples & compute stat of interest (corr in this case)	<code>B = 1000; n = nrow(data);</code> <code>theta.b = numeric(B);</code> <code>for (b in 1:B) {</code> <code>index = sample(1:n, size=n, replace=TRUE)</code> <code>bsample1 = var1[index]</code> <code>bsample2 = var2[index]</code> <code>theta.b[b] = corr(bsample1, bsample2)</code> <code>}</code>
Bootstrap Estimation of Standard Error	Compute sd	<code>sd(theta.b)</code>
Bootstrap Estimation of Bias	Compute bias	<code>bias = mean(theta.b) - theta.hat</code>
Alternate method for SE and Bias	Using Library boot	<code>bcor = function (data, b) {</code> <code>return(corr(data[b,1],data[b,2])) }</code> <code>library(boot)</code> <code>boot.corr = boot(data, statistic=bcor, R=5000)</code> <code>#gives original value from sample, bias and standard error</code>
Basic bootstrap Confidence interval	Compute lower and upper quantile	<code>alfa = 0.05</code> <code>low = quantile(theta.b,alfa/2)</code> <code>high = quantile(theta.b,1-alfa/2)</code>
	Compute CI	<code>lower = 2*theta.hat-high</code> <code>upper = 2*theta.hat-low</code>
Percentile bootstrap confidence interval	Compute lower and upper quantile	<code>alfa = 0.05</code> <code>low = quantile(theta.b,alfa/2)</code> <code>high = quantile(theta.b,1-alfa/2)</code>
Alternate method for CI	Using Library boot	<code>bcor = function (data, b) {</code> <code>return(corr(data[b,1],data[b,2])) }</code> <code>library(boot)</code> <code>boot.corr = boot(data, statistic=bcor, R=5000)</code> <code>boot.ci(boot.corr,type=c("basic","norm","perc"))</code> <code>#gives normal, basic and percentile CI</code> <code>#must generate boot.corr object first</code>
Numerical methods in R		
Finding Root	One-dimension *uniroot only finds 1 root in region *function must have opposite signs in two end points	<code>a = 0.5; n = 20;</code> <code>fy = function(y) {a^2+y^2+2*a*y/(n-1)-(n-2)}</code> <code>rt = uniroot(fy,lower=0,upper=100)</code> <code>root = rt\$root #to get root</code> <code>value = rt\$f.root #to get the value of fn</code>
	Finding roots *polyroot finds multiple roots, including complex roots	<code>#input the coefficient of polynomial</code> <code>a0 = a^2 - (n-2); a1 = 2*a/(n-1); a2 = 1;</code> <code>polyroot(c(a0,a1,a2))</code> <code>#a_i is the coefficient of xⁱ</code>
Integration	One variable integral	<code>fcn = function(y,mu,sigma) {</code> <code>abs(y)*exp(-(y-mu)^2/(2*sigma^2)) /</code> <code>sqrt(2*pi)/sigma }</code> <code>integrate(fcn,lower=-Inf, upper=Inf, mu=0,</code> <code>sigma=1)</code>
Optimiation	One-dimentional	<code>f = function(x) {log(x+log(x))/log(1+x)}</code> <code>optimize(f,lower=4,upper=8,maximum=TRUE)</code>

Function; Variable; Reserved Keywords		R
	Two-dimensional *”optim” performs minimization by default	<pre> LL = function(theta, sx, slogx, n) { r = theta[1]; lambda = theta[2] loglik = n*r*log(lambda) + (r-1)*slogx - lambda*sx - n*log(gamma(r)) return(loglik)} #or return(-loglik) if lazy to use fnscale n = 200; r = 5; lambda = 2; theta.start = c(1,1) x = rgamma(n, shape=r, rate=lambda) out = optim(par=theta.start, fn=LL, sx=sum(x), slogx=sum(log(x)), n=n, control=list(fnscale=-1)) r = out\$par[1]; lambda = out\$par[2] result = out\$value </pre>
Regression Analysis		
Pearson Correlation Coefficient		<pre> cor(dataframe, method="person") #test if true correlation is 0 cor.test(var1, var2) </pre>
Simple regression		<pre> modell = lm(y~x) #to see sum of squared regression and error anova(modell) #to see R-squared and linear model, F-test summary(modell) </pre>
Testing assumptions - Linear relationship: scatter plot		<pre> modell = lm(y~x) par(mfrow=c(2,2)) #scatter plot plot(y~x); abline(modell); title("Scatter plot and Regression Line") </pre>
- Normality and independence observe if there's any trend in residuals, it shouldn't have		<pre> #residual plot rs = modell\$resid; fv = modell\$fitted; plot(rs~x); abline(h=0) #normal QQ plot qqnorm(rs); qqline(rs) par(mfrow=c(1,1)) </pre>
- Prediction	case1: data in set	<code>modell\$coeff[1]+value*modell\$coeff[2]</code>
	case2: data not in	<pre> predict(modell, data.frame(x=value)) #or the previous method </pre>
- Transformation	Investigate rs	<code>plot(x, y); lines(smooth.spline(x, y))</code>
	quadratic	<pre> x2 = x^2; model2 = lm(y~x+x2) anova(model2); summary(model2) </pre>
	log	<pre> lx = log(x); model3 = lm(y~lx) anova(model3); summary(model3) </pre>

Function; Variable; Reserved Keywords			done by: Ling 1 Dec 19
		SAS	R
Simulation			
Generate random numbers	Linear congruential generators	<pre>data try1; seed = 1234; do i = 1 to 10; x = ranuni(seed); output; end; keep x;</pre>	<pre>ran = NULL for (i in 1:n){ x1 = (a*x0+c)%m x0 = x1 ran = c(ran,x1/m) }</pre>
	Uniform random numbers	<pre>data unif; seed = 1234; call streaminit(seed); n=100; a=0; b=10; do i = 1 to n; x = a + (b-a) * rand("UNIFORM"); output; end; keep x;</pre>	<pre>#uniform with (0,1) x = runif(1000) #uniform with (a,b) set.seed(1234) x = runif(n,a,b)</pre>
	Normal(μ, σ^2)	<pre>data norm; seed = 1234; call streaminit(seed); n=100; mu=0; sigma=1; do i = 1 to n; x = rand("NORMAL", mu, sigma); output; end; keep x;</pre>	<pre>x = rnorm(n,mean=mu,sd=sigma)</pre>
	Exponential(λ)	<pre>data expno; seed = 1234; call streaminit(seed); n=100; lambda=5; do i = 1 to n; x = lambda * rand('EXPONENTIAL'); output; end; keep x;</pre>	<pre>x = rexp(n,mean=lambda)</pre>
	Gamma(α, β)	<pre>data gammano; seed = 1234; call streaminit(seed); n=100; alpha=1; beta=2; do i = 1 to n; x = beta * rand('GAMMA', alpha); output; end; keep x;</pre>	<pre>x = rgamma(n,shape=alpha,scale=beta)</pre>
	Chi-square(p)	<pre>data chisqno; seed = 1234; call streaminit(seed); n=100; df=10; do i = 1 to n; x = rand('CHISQUARE', df); output; end; keep x;</pre>	<pre>x = rchisq(n,df=p)</pre>

	Beta(α, β)	<pre> data betano; seed = 1234; call streaminit(seed); n=100; alpha=2; beta=3; do i = 1 to n; x = rand('BETA', alpha, beta); output; end; keep x; </pre>	<pre> x = rbeta(n, shape1=a, shape2=b) </pre>
	t(k)	<pre> data tno; seed = 1234; call streaminit(seed); n=100; df=5; do i = 1 to n; x = rand('T', df); output; end; keep x; </pre>	<pre> x = rt(n, df=k) </pre>
	F(m,n)	<pre> data fno; seed = 1234; call streaminit(seed); n=100; df1=5; df2=10; do i = 1 to n; x = rand('F', df1, df2); output; end; keep x; </pre>	<pre> x = rf(n, df1 = n1, df2 = n2) </pre>
	Binomial(n,p)	<pre> data binomno; seed = 1234; call streaminit(seed); ns=100; n=10; p=0.3; do i = 1 to ns; x = rand('BINOMIAL', p, n); output; end; keep x; </pre>	<pre> x = rbinom(n, size, prob = p) </pre>
	Poisson(λ)	<pre> data poisno; seed = 1234; call streaminit(seed); n=100; lambda=3; do i = 1 to n; x = rand('POISSON', lambda); output; end; keep x; </pre>	<pre> x = rpois(n, lambda) </pre>
	Hypergeo(n,N,S)	<pre> data hypergeono; seed = 1234; call streaminit(seed); ns=100; popsize=3; numbsucc=5; samplesize=3; do i = 1 to ns; x = rand('HYPERGEOMETRIC', popsize, numbsucc, samplesize); output; end; keep x; </pre>	<pre> x = rhyper(nn, m=S, n=N, k=n) </pre>

	NBinom(r,p)	<pre> data negbinomno; seed = 1234; call streaminit(seed); n=100; p=0.3; k=5; do i = 1 to n; x = rand('NEGBINOMIAL', p, k); output; end; keep x; </pre>	<pre> x = rnbinom(n, size=r, prob=p) </pre>
Simulation Study			
Conditon of interest - Compare 3 estimators for location μ through a simulation study - 3 estimators are: sample mean, sample median, sample 10% trimmed mean - Underlying condition: N(1,1) - Sample size: 15 - Simulation size: 1000			
Size of t-test α level	Generate sample set	<pre> data simu; seed = 123; ns=1000; n=15; mu=0; sigma=1; call streaminit(seed); do mcrep = 1 to ns; do i = 1 to n; value = rand("NORMAL",mu,sigma); output; end; keep mcrep value; end; </pre>	
	Calculate t test	<pre> proc sort data=simu; by mcrep; proc univariate data=simu noprint mu0=0; by mcrep; var value; output out=outtest probt = p; </pre>	
	Calculate rejection	<pre> data outtest1; set outtest; reject = (p<0.05); </pre>	
	Calculate rejection rate	<pre> proc means data=outtest1 noprint; var reject; output out = results mean=rejrate; proc print data=results; var _freq_ rejrate; </pre>	

done by: Ling 1 Dec 19			
Describing numerical data			
Numerical descriptive measures	- Location	- Mean - Median - Mode	
	- Variability	- Variance or standard deviation - Range - Interquartile range - Coefficient of Variation	
	- Other measures unbiased estimator of kurtosis: $\frac{(n-1)}{(n-2)(n-3)} \left(\frac{(n+1)m_4}{m_2^2} - 3 \left(\frac{n}{n-1} \right) \right)$	- Minimum and Maximum - First quartile and third quartile - Percentiles - Skewness [measure the symmetric] - Kurtosis [measure the tail]	unbiased estimator of skewness: $\frac{\sqrt{n(n-1)}}{n-2} \left(\frac{m_3}{m_2^{3/2}} \right)$
Graphical methods	- Histogram		
	- Boxplot, Stem and leaf plot		
	- Scatter plot for bivariate data		
Robust Location estimator			
Robust location estimators $\sum \psi(y_i - \mu) = 0$	- Trimmed mean (top-left)	mean: $\psi(x) = x$ median: $\psi(x) = \text{sign}(x)$ trimmed: $\psi(x) = x \mid x < c, \text{ else } 0$ winsorized: $\psi(x) = -c, x < -c \text{ e.g. } g$ turkey: $\psi(x) = x \left[1 - \left(\frac{x}{R} \right)^2 \right]^2 \text{ for } X < R$ 	
	- Winsorized mean		
	- Huber's M-estimators (top-right)		
	- Tukey's bisquare estimator (bottom-left)		
	- Hampel's M-estimator (bottom-right)		
Robust measures of scale parameter	- Interquartile Range (IQR)	Normal: $\hat{\sigma} = IQR / 1.34898$	
	- Median Absolute Deviation (MAD) (median of the diff in median)	$MAD = \text{median}_i(y_i - \text{median}_j(y_j))$ Normal: $\hat{\sigma} = 1.4826 \times MAD$	
	- Gini's mean difference	$G = \frac{1}{\binom{n}{2}} \sum_{i < j} y_i - y_j $ $\sigma \approx \frac{G}{2} \sqrt{\pi}$ e.g. $G = (a_1 + \dots + a_n) / n$	
Categorical data			
	- Frequency tables		
	- Contingency tables		
	- Chi-square tests (single tail test)	$H_0: \mu_1 \text{ and } \mu_2 \text{ are independent}$ $H_1: \mu_1 \text{ and } \mu_2 \text{ are not independent}$	
	- Charts		
Paired Data	Same subjects under two different conditions		
	- McNemar's test	$H_0: \text{before} = \text{after}$ $H_1: \text{before} \neq \text{after}$	
	- Bar Chart		
Numerical data			
One Sample Tests (Hypothesis testing)			
Location	Parametric test - One sample t-test (normal population)	$\mu_0 = 0$ $\mu_0 \neq 0$	

	Nonparametric tests - Sign test - Sign rank test	$\mu_0 = 0$ $\mu_0 \neq 0$
Two comparing two groups (Hypothesis testing)		
Independent samples	Two-sample t-test (normal population)	$\mu_0 = \mu_1$ $\mu_0 \neq \mu_1$
	Wilcoxon rank-sum test (non-parametric)	
Related samples	Paired t-test (normal population)	$\mu_{before} = \mu_{after}$ $\mu_{before} \neq \mu_{after}$
	Sign Test or Wilcoxon signed rank test (non-parametric)	
3 or more independent groups (One-way analysis of variance)		
Parametric method	One-way analysis of variance - ANOVA	$\mu_x = \mu_y = \mu_z$ $\mu_x \neq \mu_y \neq \mu_z$
	Assumptions of ANOVA (based on F-test) & Model Checking	<ol style="list-style-type: none"> 1. Random sample 2. Equal variances for all the groups <ol style="list-style-type: none"> a. Bartlett test or Levene's test 3. Independence of errors <ol style="list-style-type: none"> a. Residual plots <ol style="list-style-type: none"> i. QQplot on residuals ii. Plot residuals against the groups 4. Normal distribution of errors <ol style="list-style-type: none"> a. Normality test on residuals 5. Additivity of treatment effects
	Multiple comparison tests - Least significant difference (LSD)	LSD for $\mu_i - \mu_j$ $t_{\frac{\alpha}{2}, N-k} \sqrt{MSE \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$ We conclude μ_i is different from μ_j if $ \bar{X}_i - \bar{X}_j > LSD$
	- Contrast method	A contrast of means is a linear combination of the means such that the coefficients of the means sum up to zero e.g. $C_1 = \mu_2 - \mu_3$, $C_2 = 2\mu_1 - \mu_2 - \mu_3$ $H_0: C_1 = 0 \text{ and } C_2 = 0$
	Others - Duncan's multiple-range test - Student-Newman-Keul's multiple-range test - Scheffe's multiple-comparison procedure	
Nonparametric method	Kruskal-Wallis test	
Simulation		
Usage	- Checking distribution theory	<u>Theory:</u> Take a sample of size 4 from normal distribution with mean 0 $T = \frac{\bar{X}}{S/\sqrt{4}} \sim t(3)$ <u>Simulation:</u> Generate 1000 random samples with size 4. Compute statistic T and construct a histogram. Superimpose the t(3) curve and compare

	- Comparing estimators	<u>Question:</u> Comparing robustness of estimators in different conditions (underling distributions) Consider $MSE = Bias^2 + Variance$ <u>Simulation:</u> Generate random samples. Compute MSE for each estimators.
	- Buffon's needle experiment	<u>Question:</u> Compute a question of interest through a function. Consider the probability that needle intersects a line $theory: \frac{2/l}{\pi d}$ <u>Simulation:</u> Simulate the experiement by N times and count the chances of success n/N give the estimate of the probability of success
Generate numbers		
Random number	- Congruential generators	$def: X_i = (aX_{i-1} + c) \bmod M$ Uniform random numbers: $U_i = X_i/M$
Uniform numbers	- Uniform random numbers	
Non-uniform random numbers	- Inversion method	If X has a continuous distribution function F(x), where $F(x) = \Pr(X \leq x)$, then $F(x) \sim Uniform(0,1)$ Generate U from Uniform(0,1) and set $X = F^{-1}(U)$
	- Exponential distribution	$Let U = F(X) = 1 - e^{-x/\lambda}$ $Then X = F^{-1}(U) = -\lambda \log(1 - U)$
	- Weibull distribution	$F(x) = 1 - \exp(-x^\beta), D = (0, \infty)$ $X = (-\log(1 - U))^{1/\beta}$
	- Cauchy distribution	$F(x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1}\left(\frac{x - \mu}{\sigma}\right)$ $X = \sigma \tan[\pi(U - 0.5)] + \mu$
Generate Normal random variable	- Box-Muller Algorithm	Generate U_1 and U_2 from uniform(0,1) set $\theta = 2\pi U_1$; $R = (-2\log U_2)^{1/2}$ set $X = R \cos(\theta)$ and $Y = R \sin(\theta)$ Now we have X and Y, two independent standard normal variables
Random variable from other random variables	- Cauchy distributions	$Y, Z \text{ independent and } \sim N(0,1)$ $X = \frac{Y}{Z} \sim Cauchy(0,1)$ $\text{If } Y \sim N(\mu, \sigma^2) \text{ and } Z \sim N(0,1)$ $X \sim Cauchy(\mu, \sigma^2)$
	- Chi-square distribution	$\text{If } Y \sim N(0,1)$ $X = Y^2 \sim \chi^2(1)$ $\text{If } Y_1 \dots Y_n \text{ are i. i. d and } X = Y_1^2 + \dots + Y_n^2$ $X \sim \chi^2(n)$
	- Student's t-distribution	$\text{If } Y \sim N(0,1) \text{ and } Z \sim \chi^2(p), \text{ then}$ $X = \frac{Y}{\sqrt{\frac{Z}{p}}} \sim t(p)$
	- F distribution	$Y \sim \chi^2(m) \text{ and } Z \sim \chi^2(n)$ $X = \frac{Y/m}{Z/n} \sim F(m, n)$
Functions	Other functions to generate the random no.	https://stat.ethz.ch/R-manual/R-devel/library/stats/html/Distributions.html http://support.sas.com/documentation/cdl/en/lrdict/64316/HTML/default/viewer.htm#a001466748.htm
Generate random numbers	Uniform(a,b)	$f(x) = \frac{1}{b-a}, a < x < b$
	Normal(μ, σ^2)	$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), -\infty < x < \infty$

	Exponential(λ)	$f(x) = \frac{1}{\lambda} \exp\left(-\frac{x}{\lambda}\right), x > 0$
	Gamma(α, β)	$f(x) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} \exp\left(-\frac{x}{\beta}\right), x > 0$
	Chi-square(p)	$f(x) = \frac{1}{2^{\frac{p}{2}} \Gamma(\frac{p}{2})} x^{\frac{p}{2}-1} \exp\left(-\frac{x}{2}\right), x > 0$
	Beta(α, β)	$f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, \text{ for } 0 < x < 1$
	t(k)	$f(x) = \frac{\Gamma(\frac{k+1}{2})}{\Gamma(\frac{k}{2})} \frac{1}{\sqrt{k\pi}} \frac{1}{\left(1 + \frac{x^2}{k}\right)^{\frac{k+1}{2}}}, \text{ for } -\infty < x < \infty$
	F(m,n)	$f(x) = \frac{\Gamma\left(\frac{n_1 + n_2}{2}\right)}{\Gamma\left(\frac{n_1}{2}\right)\Gamma\left(\frac{n_2}{2}\right)} \left(\frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \frac{x^{\frac{n_1}{2}-1}}{\left(1 + \frac{n_1}{n_2}x\right)^{\frac{n_1+n_2}{2}}}, 0 < x < \infty$
	Binomial(n,p)	$f(x) = \binom{n}{k} p^x (1-p)^{n-x}, \text{ for } x = 0, 1, \dots, n$
	Poisson(λ)	$f(x) = \frac{e^{-\lambda} \lambda^x}{x!}, x = 0, 1, 2, \dots$
	Hypergeo(n,N,S)	$f(x) = \frac{\binom{S}{x} \binom{N-S}{n-x}}{\binom{N}{n}}, \text{ for } x = 0, 1, \dots, \min(n, S)$
	NBinom(r,p)	$f(x) = \binom{r+x-1}{x} p^r (1-p)^x, \text{ for } x = 0, 1, 2, \dots$
Simulation Studies in Stats		
Rationale (in statistics)	Properties of estimators	<ul style="list-style-type: none"> - Bias - Consistent - Sampling variance - Comparison with competing estimators on bias, precision etc - Does confidence interval achieve the advertised nominal level of coverage (e.g.95%)?
	Properties of hypothesis testing procedures	<ul style="list-style-type: none"> - Does hypothesis testing procedure attain the advertised level of significance or size (e.g. $\alpha = 0.05$)? - What power is possible against different alternatives to the null hypothesis? Do different test procedures deliver different power?
Properties of estimators	Monte Carlo simulation approximation	<ul style="list-style-type: none"> - Generate S independent data sets under the conditions of interest - Compute the numerical value of the estimators - If S large enough, summart statistics should a good approximations to the true properties of the estimator <u>under the conditions of interest</u>
	Checking the coverage probability of confidence interval	$\left(\bar{Y} - t_{0.025, n-1} \frac{S}{\sqrt{n}}, \bar{Y} + t_{0.025, n-1} \frac{S}{\sqrt{n}}\right)$
Properties of hypothesis tests	Testing size/ level of test	Generate data under $H_0: \mu = \mu_0$ Calculate the proportion of rejections of H_0 Approximate the true probability of rejecting H_0 when it is true. Proportion should be $\approx \alpha$
	Evaluate power	Generate data under $H_1: \mu \neq \mu_0$ Calculate the proportion of rejections of H_0 Approximate the true probability of rejecting H_0 when it is false.
Bootstrap Method		
Bootstrap distribution	Empirital distribution	$f_n(x) = \begin{cases} 1/n, & \text{for } x = x_1 \dots x_n \\ 0, & \text{otherwise} \end{cases}$
	Empirical cumulative distribution	$F_n(t) = \frac{\# \text{ of } x's < t}{n}$

Method	Estimate distribution function $F_{\hat{\theta}}(\cdot)$	Generate bootstrap replicate $x^{*(b)} = (x_1^*, \dots, x_n^*)$ into $\hat{\theta}^{*(b)} = (\hat{\theta}_1^*, \dots, \hat{\theta}_n^*)$ Calculate bootstrap estimate of $F_{\hat{\theta}}(\cdot)$
Statistics	Standard Error	$\widehat{se}(\hat{\theta}^*) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}^{*(b)} - \bar{\hat{\theta}}^*)^2}$ <p>where</p> $\bar{\hat{\theta}}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^{*(b)}$
	Bias	$bias(\hat{\theta}) = E(\hat{\theta} - \theta) = E(\hat{\theta}) - \theta$ $\widehat{bias}(\hat{\theta}) = \bar{\hat{\theta}}^* - \hat{\theta}$
	Confidence interval	<p>Basic</p> $(2\hat{\theta} - \hat{\theta}_{1-\alpha/2}^*, 2\hat{\theta} - \hat{\theta}_{\alpha/2}^*)$ <p>Percentile</p> $(\hat{\theta}_{\alpha/2}^*, \hat{\theta}_{1-\alpha/2}^*)$ <p>where $\hat{\theta}_{\alpha}^*$ is the α sample quantile from the empirical cdf of the replicates $\hat{\theta}^*$</p>
Numerical methods in R		
Root-finding	One-dimension	e.g. solve for y in equation $a^2 + y^2 + \frac{2ay}{n-1} = n-2, a = 0.5, n = 20$
Integration	One-variable	e.g. solve for the integral $\int \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$
Maximum Likelihood Method	Concept	A method to estimate the parameters of a distribution, the parameter should maximize the likelihood function $\max L(\theta) = f(x_1, x_2, \dots, x_n \theta)$ Assuming independent and identically distributed random values $f(x_1, x_2, \dots, x_n \theta) = \prod_{i=1}^n f(x_i \theta)$ <p>Since population parameter θ is unknown, we estimate with $\hat{\theta}$ In reality the product function is difficult to differentiate, so we take log of the function</p>
	Method1: finding θ which maximizes $L(\theta)$ or $l(\theta)$	Newton's method to find the value which maximizes the function
	Method2: finding the first order condition	$\frac{d}{d\theta} L(\theta) = 0 \text{ or } \frac{d}{d\theta} l(\theta) = 0$
Optimization	One-dimension	e.g. Maximize the function $f(x) = \frac{\log(1 + \log(x))}{\log(1+x)}, \text{ wrt } x \text{ and } x > 0$
	Two-dimension	e.g. Find the maximum likelihood estimator of r and λ $\max L(r, \lambda) = \frac{(\lambda)^{nr}}{\Gamma(r)^n} \prod_{i=1}^n x_i^{r-1} \exp\left(-\lambda \sum_{i=1}^n x_i\right), x_i \geq 0$ or consider the log-likelihood function $\max l(r, \lambda) = nr \log \lambda - n \log \Gamma(r) + (r-1) \sum_{i=1}^n \log x_i - \lambda \sum_{i=1}^n x_i$
Regression Analysis		
Correlation	Pearson Correlation Coefficient - Indicates the strength of linear relationship between two variables	$\rho = \frac{Cov(X, Y)}{\sqrt{Var(x)}\sqrt{Var(Y)}}$ $\rho = \frac{E[(X - \mu_x)(Y - \mu_y)]}{\sqrt{E[(X - \mu_x)^2]} \sqrt{E[(Y - \mu_y)^2]}}$

Sample regression	Linear model establish a relationship between two variables	$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ where $\varepsilon_i \sim N(0, \sigma^2)$ independently $E(Y_i x_i) = \beta_0 + \beta_1 x_i$ $Var(Y_i x_i) = Var(\varepsilon_i) = \sigma^2$
	Assumptions	<ul style="list-style-type: none"> - Y_i's are independent - Y_i's follow normal distribution with given parameters Study residuals <ul style="list-style-type: none"> - Y_i's have the same variances Study residuals <ul style="list-style-type: none"> - Linear relationship exist between Y and X Test using scatter plot
	Prediction	Case1: values in the dataset Case2: values not within the dataset
	Transformation	<ul style="list-style-type: none"> - Quadratic term to the regression line - Log transformation
Bonus		
Machine Learning	Extra	https://www.analyticsvidhya.com/blog/2017/09/common-machine-learning-algorithms/

Data Transform Analyze Graphs Utilities

- Define Variable Properties...
- Set Measurement Level for Unknown...
- Copy Data Properties...
- New Custom Attribute...
- Define date and time...
- Define Multiple Response Sets...

Validation

- Identify Duplicate Cases...
- Identify Unusual Cases...
- Compare Datasets...

- Sort Cases...
- Sort Variables...
- Transpose...
- Adjust String Widths Across Files

Merge Files

- Restructure...
- Rake Weights...
- Propensity Score Matching...
- Case Control Matching...

- Aggregate...

Orthogonal Design

- Split into Files

- Copy Dataset

- Split File...

- Select Cases...

- Weight Cases...

Transform Analyze Graphs Utilities Extension

- Compute Variable...
- Programmability Transformation...
- Count Values within Cases...
- Shift Values...

- Recode into Same Variables...
- Recode into Different Variables...
- Automatic Recode...

- Create Dummy Variables

- Visual Binning...

- Optimal Binning...

Prepare Data for Modeling

- Rank Cases...

- Date and Time Wizard...




- Create Time Series...

- Replace Missing Values...

- Random Number Generators...

- Run Pending Transforms

Ctrl+G

Analyze	Graphs	Utilities	Extensions
Reports			►
Descriptive Statistics			►
Bayesian Statistics			►
Tables			►
Compare Means			►
General Linear Model			►
Generalized Linear Models			►
Mixed Models			►
Correlate			►
Regression			►
Loglinear			►
Neural Networks			►
Classify			►
Dimension Reduction			►
Scale			►
Nonparametric Tests			►
Forecasting			►
Survival			►
Multiple Response			►
 Missing Value Analysis...			
Multiple Imputation			►
Complex Samples			►
 Simulation...			
Quality Control			►
 ROC Curve...			
Spatial and Temporal Modeling...			►
Direct Marketing			►