

STATISTICS with R

ANALYSIS OF ONE VARIABLE

Continuous variable	<p><i>x continuous variable</i></p> <ul style="list-style-type: none"> ○ Summary statistics <i>summary(x) # most important summary statistics</i> <i>min(x) # minimum</i> <i>max(x) # maximum</i> <i>mean(x) # mean, average</i> <i>median(x) # median</i> <i>sd(x) # standard deviation</i> <i>IQR(x) # interquartile rang</i> <i>quantile(,) # Ex. 95% percentile: quantile(x, 0.95)</i> ○ Dot plot <i>plot(x)</i> ○ Histogram <i>hist(x)</i> ○ Box plot or Box-and-whisker plot <i>boxplot(x)</i> ○ Density function <i>plot(density(x))</i> ○ Empirical cumulative distribution <i>plot(ecdf(x))</i>
Categorical variable	<p><i>x categorical variable</i></p> <ul style="list-style-type: none"> ○ Frequency table <i>table(x)</i> <i>prop.table(table(x)) # Table of relative frequencies</i> <i>100*prop.table(table(x)) # Table of percentages</i> ○ Bar plot <i>barplot(table(x))</i> ○ Pie chart <i>pie(table(x))</i>

RELATION BETWEEN TWO VARIABLES

	Relation between two variables
Continuous & continuous	<p><i>x and y continuous variables</i></p> <ul style="list-style-type: none"> ○ Correlation coefficient <code>cor(x,y)</code> # Pearson correlation coefficient <code>cor(x,y, method="spearman")</code> # Spearman correlation coefficient <code>cor(M, use="pairwise.complete.obs")</code> # M is a matrix ○ Regression line equation <code>lm(y~x)</code> ○ Scatter plot and regression line <code>plot(x,y)</code> # independent before dependent (x,y) <code>abline(lm(y~x))</code> # dependent before independent (y,x)
Continuous & categorical	<p><i>y continuous, x categorical</i></p> <ul style="list-style-type: none"> ○ Numerical summaries of the continuous variable by each category of the categorical variable <code>tapply(<continuous>, <categorical>, <function>)</code> # Example: <code>tapply(y, x, mean)</code> # mean of y for each category of x <code>tapply(y, x, summary)</code> # summary of y for each category of x ○ Multiple box plot <code>boxplot(<continuous> ~<categorical>)</code> # Example: <code>boxplot(y~x)</code>
Categorical & categorical	<p><i>x and y categorical variables</i></p> <ul style="list-style-type: none"> ○ 2 by 2 table / Contingency table <code>table(x,y)</code> # absolute frequencies <code>prop.table(table(x,y))</code> # total proportions <code>prop.table(table(x,y),1)</code> # row proportions <code>prop.table(table(x,y),2)</code> # column proportions <code>100*prop.table(table(x,y),1)</code> # row percentages ○ Bar plot <code>barplot(table(x,y))</code> <code>barplot(prop.table(table(x,y)))</code>

RANDOM VARIABLES WITH R

$f(x) \text{ or } P(X = x)$ $P(X \leq x)$ $P(X \leq q) = \alpha$

Table 3.2: Built-in-functions for random variables used in this chapter.

Distribution	parameters	density	distribution	quantiles	random sampling
Bin	n, p	<code>dbinom(x, n, p)</code>	<code>pbinom(x, n, p)</code>	<code>qbinom(α, n, p)</code>	<code>rbinom(10, n, p)</code>
Normal	μ, σ	<code>dnorm(x, μ, σ)</code>	<code>pnorm(x, μ, σ)</code>	<code>qnorm(α, μ, σ)</code>	<code>rnorm(10, μ, σ)</code>
Chi-squared	m	<code>dchisq(x, m)</code>	<code>pchisq(x, m)</code>	<code>qchisq(α, m)</code>	<code>rchisq(10, m)</code>
T	m	<code>dt(x, m)</code>	<code>pt(x, m)</code>	<code>qt(α, m)</code>	<code>rt(10, m)</code>
F	m, n	<code>df(x, m, n)</code>	<code>pf(x, m, n)</code>	<code>qf(α, m, n)</code>	<code>rf(10, m, n)</code>

- **Other distributions:**

Geometric: `dgeom()`

Negative Binomial: `dnbinom()`

Poisson: `dpois()`

Hipergeometric: `dhyper()`

Exponential: `dexp()`

- **Examples Binomial distribution**

X Binomial with parameters $n = 8$ i $p = 0.35$

$P(X = 4)$: `dbinom(4, 8, 0.35)`

$P(X \leq 4)$: `pbinom(4, 8, 0.35)`

95% Percentile: `qbinom(0.95, 8, 0.35)`

Random sample of 25 values of X : `rbinom(25, 8, 0.35)`

- **Examples Normal distribution**

X Normal of parameters $\mu = 10$ i $\sigma = 3$

$P(X \leq 15)$: `pnorm(15, 10, 3)`

$P(X > 20)$: `1-pnorm(20, 10, 3)`

$P(12 \leq X \leq 20)$: `pnorm(20, 10, 3) - pnorm(12, 10, 3)`

95% Percentile: `qnorm(0.95, 10, 3)`

Random sample of 25 values of X : `rnorm(25, 10, 3)`

STATISTICAL TESTS WITH R

<i>y continuous variable</i> <i>x categorical variable</i>	Normality Test: Shapiro-Wilk H0: Data follow a normal distribution H1: Data do not follow a normal distribution <i>tapply(<continuous>,<categorical>,function(x) shapiro.test(x))</i>	
	If Shapiro p-value > alpha Data follow a normal distribution	If Shapiro p-value < alpha Data DO NOT follow a normal distribution
Test for the mean H0: mean=prespecified value H1: mean≠ prespecified value	T-test t for one sample <i>t.test(y, mu=value)</i>	Wilcoxon test for one sample <i>wilcox.test(y, mu=value)</i>
Test for the equality of two means H0: mean1=mean2 H1: mean1≠ mean2	T-test for independent samples (previously, you should test for the equality of variances) <i>t.test(y~x, var.equal=T) # if variances are equal</i> <i>t.test(y~x,var.equal=F) # if variances are different</i>	Wilcoxon test for independent samples (also known as Wilcoxon–Mann–Whitney test) <i>wilcox.test(y~x)</i>
Test for the equality of two means with paired samples H0: mean1=mean2 H1: mean1≠ mean2	T-test for paired samples <i>d<-y1-y2</i> <i>t.test(d,mu=0)</i>	Wilcoxon test for paired samples <i>wilcox.test(y1,y2,paired=TRUE)</i>
Test for the equality of more than two means H0: mean1 = mean2 = ... = meank H1: at least one of the means is different	one-factor ANOVA (Requires normality and homoscedasticity) <i>aov(y~x)</i> Normality: <i>shapiro.test(residuals(lm(y~x)))</i> Post-hoc analysis: <i>TukeyHSD(aov)</i> Robust ANOVA (if homoscedasticity is not fulfilled): <i>oneway.test(y~x)</i> two-factor ANOVA <i>aov(y~x1*x2)</i>	Kruskal-Wallis test <i>kruskal.test(y~x)</i>
Test for the equality of two variances H0: variance1= variance2 H1: variance1≠ variance2	F test for the equality of variances <i>var.test(y~x)</i>	
Test for the equality of several variances H0: var1 = var2 = ... = vark H1: at least one of the means is different	Homoscedasticity test <i>install.packages("lmtest")</i> <i>library(lmtest)</i> <i>bptest(lm(y ~ x),studentize = F)</i>	

Test for one proportion H0: $p = \text{prespecified value } p_0$ H1: $p \neq p_0$	Binomial test for one proportion <code>binom.test(k,n,p0)</code>
Test for equality of proportions H0: $\text{proportion1} = \text{proportion2}$ H1: $\text{proportion1} \neq \text{proportion2}$	Test for the equality of two proportions <code>prop.test(table(x1,x2))</code> # x1 i x2 are factors with 2 categories
Multinomial test $H_0: (\pi_1, \dots, \pi_m) = (p_1, \dots, p_m)$ $H_1: (\pi_1, \dots, \pi_m) \neq (p_1, \dots, p_m)$	Multinomial test for proportions <code>prop.test(x=c(n1,..., nm),p=c(p1, ..., pm))</code>
Test for independence of 2 categorical variables H0: X and Y are independent H1: X and Y are related	Chi-squared test for independence of 2 factors <code>chisq.test(table(x1,x2))</code> # x1 and x2 are categorical variables
Test for independence of 2 categorical variables with 2 categories H0: X and Y are independent H1: X and Y are related	Fisher test for independence of 2 factors (2x2 tables) <code>fisher.test(table(x1,x2))</code> # x1 and x2 are categorical variables
Test for odds ratio H0: $OR=1$ H1: $OR \neq 1$	Odds ratio test for 2 factors (2x2 tables) <code>install.packages("epitools")</code> <code>library("epitools")</code> <code>oddsratio(table(x1, x2))</code> <code>oddsratio(table, rev="c")</code> # reverse columns <code>oddsratio(table, rev="both")</code> #reverse both, columns and rows
Test for independence of two continuous variables H0: X and Y are not correlated H1: X and Y are correlated	Correlation test <code>cor.test(x,y)</code> # Pearson correlation <code>cor.test(x,y, method=c("spearman"))</code> # Spearman correlation
Outliers test H0: No outliers H1: data contain outliers	Outliers test <code>library(outliers)</code> <code>grubbs.test(x)</code>
Correction for multiple testing	Benjamini and Hochberg FDR control <code>p.adjust(p, method = "fdr", n = length(p))</code>

Regression models with R

Linear regression <i>Y continuous</i> <i>X1, X2 explanatory</i>	<pre>model<-lm(y~x1+x2, data = data) summary(model) Check normality of residuals: shapiro.test(residuals(model))</pre>
Logistic regression <i>Y binary</i> <i>X1, X2 explanatory</i>	<pre>model<-glm(y~x1+x2, data = data, family = "binomial") summary(model)</pre>
Seleccion of variables in regression <i>(step-wise regression)</i>	<pre>step(model)</pre>
Diagnostics in regression: Residuals vs predictions Plots	<pre>plot(predict(model), residuals(model)) abline(a=0, b=0)</pre>