

# STATISTICS with R

## ANALYSIS OF ONE VARIABLE

<p><b>Continuous variable</b></p>	<p><i>x continuous variable</i></p> <ul style="list-style-type: none"> <li>○ <b>Summary statistics</b>  <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> </li> <li>○ <b>Dot plot</b>  <i>plot(x)</i> </li> <li>○ <b>Histogram</b>  <i>hist(x)</i> </li> <li>○ <b>Box plot or Box-and-whisker plot</b>  <i>boxplot(x)</i> </li> <li>○ <b>Density function</b>  <i>plot(density(x))</i> </li> <li>○ <b>Empirical cumulative distribution</b>  <i>plot(ecdf(x))</i> </li> </ul>
<p><b>Categorical variable</b></p>	<p><i>x categorical variable</i></p> <ul style="list-style-type: none"> <li>○ <b>Frequency table</b>  <i>table(x)</i>  <i>prop.table(table(x)) # Table of relative frequencies</i>  <i>100*prop.table(table(x)) # Table of percentages</i> </li> <li>○ <b>Bar plot</b>  <i>barplot(table(x))</i> </li> <li>○ <b>Pie chart</b>  <i>pie(table(x))</i> </li> </ul>

# RELATION BETWEEN TWO VARIABLES

	Relation between two variables
<b>Continuous &amp; continuous</b>	<p><i>x and y continuous variables</i></p> <ul style="list-style-type: none"> <li>○ <b>Correlation coefficient</b>  <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</li> <li>○ <b>Regression line equation</b>  <code>lm(y~x)</code></li> <li>○ <b>Scatter plot and regression line</b>  <code>plot(x,y)</code> # independent before dependent (x,y)  <code>abline(lm(y~x))</code> # dependent before independent (y,x)</li> </ul>
<b>Continuous &amp; categorical</b>	<p><i>y continuous, x categorical</i></p> <ul style="list-style-type: none"> <li>○ <b>Numerical summaries of the continuous variable by each category of the categorical variable</b>  <code>tapply(&lt;continuous&gt;, &lt;categorical&gt;, &lt;function&gt; )</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</li> <li>○ <b>Multiple box plot</b>  <code>boxplot(&lt;continuous&gt; ~&lt;categorical&gt; )</code>  # Example:  <code>boxplot(y~x)</code></li> </ul>
<b>Categorical &amp; categorical</b>	<p><i>x and y categorical variables</i></p> <ul style="list-style-type: none"> <li>○ <b>2 by 2 table / Contingency table</b>  <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</li> <li>○ <b>Bar plot</b>  <code>barplot(table(x,y))</code>  <code>barplot(prop.table(table(x,y)))</code></li> </ul>

# 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(<math>x, n, p</math>)</code>	<code>pbinom(<math>x, n, p</math>)</code>	<code>qbinom(<math>\alpha, n, p</math>)</code>	<code>rbinom(10, <math>n, p</math>)</code>
Normal	$\mu, \sigma$	<code>dnorm(<math>x, \mu, \sigma</math>)</code>	<code>pnorm(<math>x, \mu, \sigma</math>)</code>	<code>qnorm(<math>\alpha, \mu, \sigma</math>)</code>	<code>rnorm(10, <math>\mu, \sigma</math>)</code>
Chi-squared	$m$	<code>dchisq(<math>x, m</math>)</code>	<code>pchisq(<math>x, m</math>)</code>	<code>qchisq(<math>\alpha, m</math>)</code>	<code>rchisq(10, <math>m</math>)</code>
T	$m$	<code>dt(<math>x, m</math>)</code>	<code>pt(<math>x, m</math>)</code>	<code>qt(<math>\alpha, m</math>)</code>	<code>rt(10, <math>m</math>)</code>
F	$m, n$	<code>df(<math>x, m, n</math>)</code>	<code>pf(<math>x, m, n</math>)</code>	<code>qf(<math>\alpha, m, n</math>)</code>	<code>rf(10, <math>m, n</math>)</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>	<b>Normality Test: Shapiro-Wilk</b> H0: Data follow a normal distribution H1: Data do not follow a normal distribution <i>shapiro.test(y)</i>	
	If Shapiro p-value >0.05 <b>Data follow a normal distribution</b>	Si Shapiro p-value <0.05 <b>Data DO NOT follow a normal distribution</b>
<b>Test for the mean</b> 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>
<b>Test for the equality of two means</b> 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>
<b>Test for the equality of two means with paired samples</b> H0: mean1=mean2 H1: mean1≠ mean2	T-test for paired samples  <i>d&lt;-y1-y2</i> <i>t.test(d, mu=0)</i>	Wilcoxon test for paired samples  <i>wilcox.test(y1,y2,paired=TRUE)</i>
<b>Test for the equality of more than two means</b> 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> 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>
<b>Test for the equality of two variances</b> H0: variance1= variance2 H1: variance1≠ variance2	F test for the equality of variances  <i>var.test(y~x)</i>	
<b>Test for the equality of several variances</b> 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>	

<b>Test for one proportion</b> H0: $p = \text{prespecified value } p_0$ H1: $p \neq p_0$	Binomial test for one proportion  <code>binom.test(k,n,p0)</code>
<b>Test for equality of proportions</b> 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
<b>Multinomial test</b> $H_0: (\pi_1, \dots, \pi_m) = (p_1, \dots, p_m)$ $H_1: (\pi_1, \dots, \pi_m) \neq (p_1, \dots, p_m)$	Multinomial test for proportions  <code>prop.test(x=c(n1,..., nm),p=c(p1, ..., pm))</code>
<b>Test for independence of 2 categorical variables</b> 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
<b>Test for independence of 2 categorical variables with 2 categories</b> 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
<b>Test for odds ratio</b> 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
<b>Test for independence of two continuous variables</b> 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
<b>Outliers test</b> H0: No outliers H1: data contain outliers	Outliers test <code>library(outliers)</code> <code>grubbs.test(x)</code>
<b>Correction for multiple testing</b>	Benjamini and Hochberg FDR control <code>p.adjust(p, method = "fdr", n = length(p))</code>