

Statistics Quick Reference Card & R Commands

by Anthony Tanbakuchi. Version 1.8.2

<http://www.tanbakuchi.com>

ANTHONY@TANBAKUCHI.COM

Get R at: <http://www.r-project.org>

R commands: **bold typewriter text**

1 Misc R

To make a vector / store data: `x=c(x1, x2, ...)`

Help: general `RSiteSearch("Search Phrase")`

Help: function `?functionName`

Get column of data from table:

`tableName$columnName`

List all variables: `ls()`

Delete all variables: `rm(list=ls())`

$$\sqrt{x} = \text{sqrt}(x)$$

$$x^n = x^{\wedge} n$$

$$n = \text{length}(x)$$

$$T = \text{table}(x)$$

2 Descriptive Statistics

2.1 NUMERICAL

Let `x=c(x1, x2, x3, ...)`

$$\text{total} = \sum_{i=1}^n x_i = \text{sum}(x)$$

$$\min = \text{min}(x)$$

$$\max = \text{max}(x)$$

six number summary : `summary(x)`

$$\mu = \frac{\sum x_i}{N} = \text{mean}(x)$$

$$\bar{x} = \frac{\sum x_i}{n} = \text{mean}(x)$$

$$\tilde{x} = P_{50} = \text{median}(x)$$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \text{sd}(x)$$

$$CV = \frac{\sigma}{\mu} = \frac{s}{\bar{x}}$$

2.2 RELATIVE STANDING

$$z = \frac{x - \mu}{\sigma} = \frac{x - \bar{x}}{s}$$

Percentiles:

$$P_k = x_i, \quad (\text{sorted } x)$$

$$k = \frac{i - 0.5}{n} \cdot 100\%$$

To find x_i given P_k , i is:

1. $L = (k/100\%)n$

2. if L is an integer: $i = L + 0.5$;
otherwise $i=L$ and round up.

2.3 VISUAL

All plots have optional arguments:

- `main=""` sets title
- `xlab="", ylab=""` sets x/y-axis label
- `type="p"` for point plot
- `type="l"` for line plot
- `type="b"` for both points and lines

Ex: `plot(x, y, type="b", main="My Plot")`

Plot Types:

`hist(x)` histogram

`stem(x)` stem & leaf

`boxplot(x)` box plot

`plot(T)` bar plot, T=table(x)

`plot(x, y)` scatter plot, x, y are ordered vectors

`plot(t, y)` time series plot, t, y are ordered vectors

`curve(expr, xmin, xmax)` plot expr involving x

2.4 ASSESSING NORMALITY

Q-Q plot: `qqnorm(x); qqline(x)`

3 Probability

Number of successes x with n possible outcomes.

(Don't double count!)

$$P(A) = \frac{x_A}{n} \quad (17)$$

$$P(\bar{A}) = 1 - P(A) \quad (18)$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B) \quad (19)$$

$$P(A \text{ or } B) = P(A) + P(B) \quad \text{if } A, B \text{ mut. excl.} \quad (20)$$

$$P(A \text{ and } B) = P(A) \cdot P(B|A) \quad (21)$$

$$P(A \text{ and } B) = P(A) \cdot P(B) \quad \text{if } A, B \text{ independent} \quad (22)$$

$$n! = n(n-1)\dots 1 = \text{factorial}(n) \quad (23)$$

$$nP_k = \frac{n!}{(n-k)!} \quad \text{Perm. no elem. alike} \quad (24)$$

$$= \frac{n!}{n_1!n_2!\dots n_k!} \quad \text{Perm. } n_1 \text{ alike, ...} \quad (25)$$

$$nC_k = \frac{n!}{(n-k)!k!} = \text{choose}(n, k) \quad (26)$$

4 Discrete Random Variables

`P(x_i)`: probability distribution

$$E = \mu = \sum x_i \cdot P(x_i)$$

$$\sigma = \sqrt{\sum (x_i - \mu)^2 \cdot P(x_i)}$$

4.1 BINOMIAL DISTRIBUTION

$$\mu = n \cdot p \quad (30)$$

$$\sigma = \sqrt{n \cdot p \cdot q} \quad (31)$$

$$P(x) = {}_n C_x p^x q^{(n-x)} = \text{dbinom}(x, n, p) \quad (32)$$

4.2 POISSON DISTRIBUTION

$$P(x) = \frac{\mu^x \cdot e^{-\mu}}{x!} = \text{dpois}(x, \mu) \quad (33)$$

5 Continuous random variables

CDF $F(x)$ gives area to the left of x , $F^{-1}(p)$ expects p is area to the left.

$f(x)$: probability density

$$E = \mu = \int_{-\infty}^{\infty} x \cdot f(x) dx \quad (34)$$

$$\sigma = \sqrt{\int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx} \quad (35)$$

$F(x)$: cumulative prob. density (CDF)

$F^{-1}(x)$: inv. cumulative prob. density

$$F(x) = \int_{-\infty}^x f(x') dx' \quad (36)$$

$$p = P(x < x') = F(x') \quad (37)$$

$$x' = F^{-1}(p) \quad (38)$$

$$p = P(x > a) = 1 - F(a) \quad (39)$$

$$p = P(a < x < b) = F(b) - F(a) \quad (40)$$

5.1 UNIFORM DISTRIBUTION

$$p = P(u < u') = F(u') \quad (41)$$

$$= \text{punif}(u', \text{min}=0, \text{max}=1) \quad (42)$$

$$u' = F^{-1}(p) = \text{qunif}(p, \text{min}=0, \text{max}=1) \quad (43)$$

5.2 NORMAL DISTRIBUTION

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}} \quad (44)$$

$$p = P(z < z') = F(z') = \text{pnorm}(z') \quad (45)$$

$$z' = F^{-1}(p) = \text{qnorm}(p) \quad (46)$$

$$p = P(x < x') = F(x') \quad (47)$$

$$= \text{pnorm}(x', \text{mean}=\mu, \text{sd}=\sigma) \quad (48)$$

$$x' = F^{-1}(p) \quad (49)$$

$$= \text{qnorm}(p, \text{mean}=\mu, \text{sd}=\sigma) \quad (50)$$

5.3 t-DISTRIBUTION

$$p = P(t < t') = F(t') = \text{pt}(t', \text{df}) \quad (51)$$

$$t' = F^{-1}(p) = \text{qt}(p, \text{df}) \quad (52)$$

5.4 χ^2 -DISTRIBUTION

$$p = P(\chi^2 < \chi^2') = F(\chi^2') \quad (53)$$

$$= \text{pchisq}(X^2', \text{df}) \quad (54)$$

$$\chi^2' = F^{-1}(p) = \text{qchisq}(p, \text{df}) \quad (55)$$

5.5 F-DISTRIBUTION

$$p = P(F < F') = F(F') \quad (56)$$

$$= \text{pf}(F', \text{df1}, \text{df2}) \quad (57)$$

$$F' = F^{-1}(p) = \text{qf}(p, \text{df1}, \text{df2}) \quad (58)$$

6 Sampling distributions

$$\mu_{\bar{x}} = \mu \quad \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \quad (59)$$

$$\mu_{\hat{p}} = p \quad \sigma_{\hat{p}} = \sqrt{\frac{pq}{n}} \quad (60)$$

7 Estimation

7.1 CONFIDENCE INTERVALS

proportion: $\hat{p} \pm E$, $E = z_{\alpha/2} \cdot \sigma_{\hat{p}}$

mean (σ known): $\bar{x} \pm E$, $E = z_{\alpha/2} \cdot \sigma_{\bar{x}}$

mean (σ unknown, use s): $\bar{x} \pm E$, $E = t_{\alpha/2} \cdot \sigma_{\bar{x}}$,

$$df = n - 1$$

$$\text{variance: } \frac{(n-1)s^2}{\chi_R^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_L^2}, \quad (61)$$

$$df = n - 1$$

$$2 \text{ proportions: } \Delta \hat{p} \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}} \quad (62)$$

$$2 \text{ means (indep): } \Delta \bar{x} \pm t_{\alpha/2} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}, \quad (63)$$

$$\text{matched pairs: } \bar{d} \pm t_{\alpha/2} \cdot \frac{s_d}{\sqrt{n}}, \quad d_i = x_i - y_i, \quad (64)$$

$$df = n - 1$$

7.2 CI CRITICAL VALUES (TWO SIDED)

$$z_{\alpha/2} = F_z^{-1}(1 - \alpha/2) = \text{qnorm}(1 - \text{alpha}/2) \quad (65)$$

$$t_{\alpha/2} = F_t^{-1}(1 - \alpha/2) = \text{qt}(1 - \text{alpha}/2, \text{df}) \quad (66)$$

$$\chi_L^2 = F_{\chi^2}^{-1}(\alpha/2) = \text{qchisq}(\text{alpha}/2, \text{df}) \quad (67)$$

$$\chi_R^2 = F_{\chi^2}^{-1}(1 - \alpha/2) = \text{qchisq}(1 - \text{alpha}/2, \text{df}) \quad (68)$$

$$\text{mean: } n = \left(\frac{z_{\alpha/2} \cdot \hat{\sigma}}{E} \right)^2 \quad (69)$$

7.3 REQUIRED SAMPLE SIZE

$$\text{proportion: } n = \hat{p} \hat{q} \left(\frac{z_{\alpha/2}}{E} \right)^2, \quad (70)$$

$(\hat{p} = \hat{q} = 0.5 \text{ if unknown})$

$$\text{mean: } n = \left(\frac{z_{\alpha/2} \cdot \hat{\sigma}}{E} \right)^2 \quad (71)$$

8 Hypothesis Tests

Test statistic and R function (when available) are listed for each.

Optional arguments for **hypothesis tests**:

alternative="two.sided" can be:
 "two.sided", "less", "greater"

conf.level=0.95 constructs a 95% confidence interval. Standard CI only when **alternative="two.sided"**.

Optional arguments for **power calculations & Type II error**:

alternative="two.sided" can be:
 "two.sided" or "one.sided"

sig.level=0.05 sets the significance level α .

8.1 1-SAMPLE PROPORTION

$H_0: p = p_0$

$$\text{prop.test}(\mathbf{x}, n, p=p_0, \text{alternative}=\text{"two.sided"}) \\ z = \frac{\hat{p} - p_0}{\sqrt{p_0 q_0 / n}} \quad (72)$$

8.2 1-SAMPLE MEAN (σ KNOWN)

$H_0: \mu = \mu_0$

$$z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} \quad (73)$$

8.3 1-SAMPLE MEAN (σ UNKNOWN)

$H_0: \mu = \mu_0$

t.test(x, mu=mu_0, alternative="two.sided")

Where **x** is a vector of sample data.

$$t = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}, \quad df = n - 1 \quad (74)$$

Required Sample size:

power.t.test(delta=h, sd = sigma, sig.level=alpha, power=1 - beta, type = "one.sample", alternative="two.sided")

8.4 2-SAMPLE PROPORTION TEST

$H_0: p_1 = p_2$ or equivalently $H_0: \Delta p = 0$

prop.test(x, n, alternative="two.sided")

where: **x=c(x1, x2)** and **n=c(n1, n2)**

$$z = \frac{\Delta \hat{p} - \Delta p_0}{\sqrt{\frac{\hat{p}\hat{q}}{n_1} + \frac{\hat{p}\hat{q}}{n_2}}}, \quad \Delta \hat{p} = \hat{p}_1 - \hat{p}_2 \quad (75)$$

$$\bar{p} = \frac{x_1 + x_2}{n_1 + n_2}, \quad \bar{q} = 1 - \bar{p} \quad (76)$$

Required Sample size:

power.prop.test(p1=p1, p2=p2, power=1 - beta, sig.level=alpha, alternative="two.sided")

8.5 2-SAMPLE MEAN TEST

$H_0: \mu_1 = \mu_2$ or equivalently $H_0: \Delta \mu = 0$

t.test(x1, x2, alternative="two.sided")

where: **x1** and **x2** are vectors of sample 1 and sample 2 data.

$$t = \frac{\Delta \bar{x} - \Delta \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad df \approx \min(n_1 - 1, n_2 - 1), \quad \Delta \bar{x} = \bar{x}_1 - \bar{x}_2 \quad (77)$$

Required Sample size:

power.t.test(delta=h, sd = sigma, sig.level=alpha, power=1 - beta, type = "two.sample", alternative="two.sided")

8.6 2-SAMPLE MATCHED PAIRS TEST

$H_0: \mu_d = 0$

t.test(x, y, paired=TRUE, alternative="two.sided")
 where: **x** and **y** are ordered vectors of sample 1 and sample 2 data.

$$t = \frac{\bar{d} - \mu_{d0}}{s_d / \sqrt{n}}, \quad d_i = x_i - y_i, \quad df = n - 1 \quad (78)$$

Required Sample size:

power.t.test(delta=h, sd = sigma, sig.level=alpha, power=1 - beta, type = "paired", alternative="two.sided")

8.7 TEST OF HOMOGENEITY, TEST OF INDEPENDENCE

$H_0: p_1 = p_2 = \dots = p_n$ (homogeneity)

$H_0: X$ and Y are independent (independence)

chisq.test(D)

Enter table: **D=data.frame(c1, c2, ...)**, where **c1, c2, ...** are column data vectors.

Or generate table: **D=table(x1, x2)**, where **x1, x2** are ordered vectors of raw categorical data.

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}, \quad df = (\text{num rows} - 1)(\text{num cols} - 1) \quad (79)$$

$$E_i = \frac{(\text{row total})(\text{column total})}{(\text{grand total})} = np_i \quad (80)$$

For 2×2 contingency tables, you can use the Fisher Exact Test:

fisher.test(D, alternative="greater")

(must specify alternative as greater)

9 Linear Regression

9.1 LINEAR CORRELATION

$H_0: \rho = 0$

cor.test(x, y)

where: **x** and **y** are ordered vectors.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y}, \quad t = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}}, \quad df = n - 2 \quad (81)$$

9.2 MODELS IN R

MODEL TYPE	EQUATION	R MODEL
linear 1 indep var	$y = b_0 + b_1 x_1$	$y \sim x_1$
...0 intercept	$y = 0 + b_1 x_1$	$y \sim 0 + x_1$
linear 2 indep vars	$y = b_0 + b_1 x_1 + b_2 x_2$	$y \sim x_1 + x_2$
...interaction	$y = b_0 + b_1 x_1 + b_2 x_2 + b_{12} x_1 x_2$	$y \sim x_1 + x_2 + x_1 x_2$
polynomial	$y = b_0 + b_1 x_1 + b_2 x_2^2$	$y \sim x_1 + I(x_2^2)$

9.3 REGRESSION

Simple linear regression steps:

1. Make sure there is a significant linear correlation.
2. **results=lm(y~x)** Linear regression of **y** on **x** vectors
3. **results** View the results
4. **plot(x, y); abline(results)** Plot regression line on data
5. **plot(x, results\$residuals)** Plot residuals

$$y = b_0 + b_1 x_1 \quad (82)$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (83)$$

$$b_0 = \bar{y} - b_1 \bar{x} \quad (84)$$

9.4 PREDICTION INTERVALS

To predict **y** when **x = 5** and show the 95% prediction interval with regression model in **results**:

predict(results, newdata=data.frame(x=5), int="pred")

10 ANOVA

10.1 ONE WAY ANOVA

1. **results=aov(depVarColName~indepVarColName, data=tableName)** Run ANOVA with data in **TableName**, factor data in **indepVarColName** column, and response data in **depVarColName** column.
2. **summary(results)** Summarize results
3. **boxplot(depVarColName~indepVarColName, data=tableName)** Boxplot of levels for factor

$$F = \frac{MS(\text{treatment})}{MS(\text{error})}, \quad df_1 = k - 1, \quad df_2 = N - k \quad (85)$$

To find required sample size and power see **power.anova.test(...)**

11 Loading and using external data and tables

11.1 LOADING EXCEL DATA

1. Export your table as a CSV file (comma separated file) from Excel.
2. Import your table into **MyTable** in R using:
MyTable=read.csv(file.choose())

11.2 LOADING AN .RDATA FILE

You can either double click on the .RData file or use the menu:

- Windows: **File**→**Load Workspace...**
- Mac: **Workspace**→**Load Workspace File...**

11.3 USING TABLES OF DATA

1. To see all the available variables type: **ls()**
2. To see what's inside a variable, type its name.
3. If the variable **tableName** is a table, you can also type **names(tableName)** to see the column names or type **head(tableName)** to see the first few rows of data.
4. To access a column of data type **tableName\$columnName**

An example demonstrating how to get the women's height data and find the mean:

```
> ls() # See what variables are defined
[1] "women" "x"
> head(women) # Look at the first few entries
   height weight
1     58    115
2     59    117
3     60    120
> names(women) # Just get the column names
[1] "height" "weight"
> women$height # Display the height data
[1] 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72
> mean(women$height) # Find the mean of the heights
[1] 65
```