

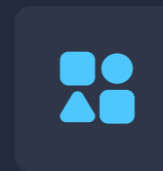


This cheat sheet provides a quick reference for essential R programming commands, helping you perform data manipulation, visualization, and statistical analysis with confidence. It covers foundational topics like installing packages and understanding R's data structures, alongside advanced tasks such as building models and applying machine learning techniques.

Each section includes concise syntax and practical examples to illustrate how R commands are used in real-world scenarios. You'll find guidance on working with vectors, lists, matrices, and data frames, performing common data wrangling tasks like filtering and summarizing, and creating visualizations such as histograms, bar plots, and boxplots. The cheat sheet also highlights R's capabilities for statistical analysis with commands like `mean`, `lm`, and `cor`.

Designed for clarity and accessibility, this resource is ideal for data analysts, statisticians, and programmers seeking to enhance their workflows in R. Whether you're exploring data, developing algorithms, or building reproducible reports, this cheat sheet ensures you can quickly apply R's powerful tools to your projects.

Table of Contents



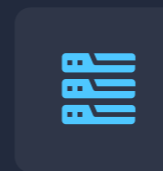
Basics

`INSTALL.PACKAGES`, `LIBRARY`,
`ASSIGNMENT (<-)`, `PRINT`, `CLASS`



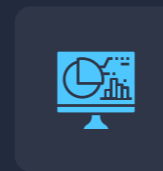
Data Structures

`C`, `LIST`, `MATRIX`, `DATA.FRAME`,
`DF$A` OR `DF`



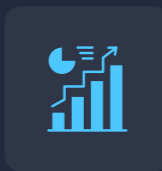
Data Manipulation

`FILTER`, `SELECT`, `MUTATE`,
`SUMMARIZE`, `ARRANGE`



Data Visualization

`PLOT`, `BARPLOT`, `HIST`, `BOXPLOT`



Statistics

`MEAN`, `MEDIAN`, `SD`, `COR`, `LM`



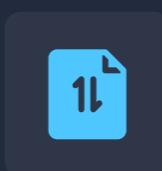
Programming

`IF`, `FOR`, `WHILE`, `FUNCTION`, `APPLY`



Machine Learning

`MATRICES`, `LINEAR MODEL`, `VISUALIZE RESIDUALS`



File I/O

`READ.CSV`, `WRITE.CSV`, `READRDS`,
`SAVERDS`, `LIST.FILES`



Basics

Syntax for	How to use	Explained
Install Package	<pre>install.packages("dplyr")</pre>	Installs the <code>dplyr</code> package.
Load Package	<pre>library(dplyr)</pre>	Loads the <code>dplyr</code> package into the current R session.
Assignment	<pre>x <- 5</pre>	Assigns value <code>5</code> to the variable <code>x</code> .
Print Output	<pre>print(x)</pre>	Prints the value of <code>x</code> to the console.
Literals and Data Types	<pre>TRUE, 125, 12.5, "Hello"</pre>	Examples of logical, integer, numeric, and character literals in R.
Extracting Numbers from Strings	<pre>library(readr) data_frame <- mutate(data_frame, column = parse_number(column))</pre>	Uses <code>parse_number</code> to extract numeric values from string columns.
Basic String Indexing	<pre>str_sub("Dataquest is awesome", 1, 9)</pre>	Extracts "Dataquest" as a substring by specifying start and end indices.

Data Structures

Syntax for	How to use	Explained
Create Vector	<pre>c(1, 2, 3)</pre>	Combines elements into a vector.
Create List	<pre>list(a=1, b="two")</pre>	Creates a list with named elements.
Create Matrix	<pre>matrix(1:6, nrow=2)</pre>	Creates a matrix with 2 rows and 3 columns.
Create Data Frame	<pre>data.frame(a=1:3, b=4:6)</pre>	Creates a data frame with columns <code>a</code> and <code>b</code> .
Access Element	<pre>df\$a df[1, 1]</pre>	Performs a logical OR operation between a column and a specific element.
Loading stringr package	<pre>library(stringr)</pre>	Loads the <code>stringr</code> library to work with strings in R.
Opening a JSON File	<pre>f <- fromJSON('filename.json')</pre>	Loads a JSON file into an R dataframe using the <code>jsonlite</code> package.
Creating a List	<pre>new_list <- list("data scientist", c(50000,40000), "programming experience")</pre>	Defines a list containing diverse data types.





Data Manipulation

Syntax for	How to use	Explained
Filter Rows	<pre>filter(df, a > 2)</pre>	Filters rows where column a is greater than 2.
Select Columns	<pre>select(df, a, b)</pre>	Selects specific columns by name.
Mutate Columns	<pre>mutate(df, c = a + b)</pre>	Adds a new column c as sum of a and b.
Summarize Data	<pre>summarize(df, avg=mean(a))</pre>	Calculates mean of column a and returns as avg.
Arrange Rows	<pre>arrange(df, desc(a))</pre>	Sorts rows by column a in descending order.
Importing Data	<pre>data <- read_csv("name_of_file_with_data.csv")</pre>	Imports dataset into R using the read_csv function from readr.
Summing Values Across Rows	<pre>df %>% mutate(new_column_name = rowSums(. [1:3]))</pre>	Sums specified columns for each row and adds as a new column.
Summing Values Across Columns	<pre>df %>% bind_rows(tibble(total = colSums(across(everything()))))</pre>	Sums specified rows for each column and adds as a new row.
Importing CSV files	<pre>dataframe <- read_csv("name_of_the_dataset.csv")</pre>	Read CSV files into R using readr's read_csv() for efficient data import.



Data Visualization

Syntax for	How to use	Explained
Creating a Basic Plot	<pre>data %>% ggplot()</pre>	Initialize a basic ggplot2 chart without specifying any aesthetics.
Creating Subplots	<pre>data %>% ggplot(aes(x = variable_1, y = variable_2)) + geom_line() + facet_wrap(~variable_3)</pre>	Plots subsets of data in separate facets.
Creating Bar Chart	<pre>data_frame %>% ggplot(aes(x = variable_1, y = variable_2)) + geom_col()</pre>	Create a bar chart using ggplot2, mapping variables to x and y axes.
Plotting multiple columns	<pre>data %>% ggplot(aes(x = variable_1)) + geom_line(aes(y = variable_2)) + geom_line(aes(y = variable_3))</pre>	Plots multiple columns on the same axes using ggplot2.
Scatterplots	<pre>ggplot(data = uber_trips, aes(x = distance, y = cost)) + geom_point()</pre>	Generate scatterplots to visualize bivariate relationships in ggplot2.
Scatterplots with Labels	<pre>ggplot(data = df, aes(x = predictor, y = response)) + geom_point() + scale_y_continuous(labels = scales::comma)</pre>	Create scatterplots with y-axis labels formatted using commas instead of scientific notation.





Data Visualization

Syntax for

How to use

Scatterplot with Comma Labels

```
ggplot(data = df, aes(x = predictor, y = response)) +
  scale_y_continuous(labels = scales::comma) + geom_point()
```

Scatterplot with Groups

```
ggplot(data = df, aes(x = predictor, y = response)) +
  geom_point() + facet_wrap(~ categorical_variable, ncol = 2)
```

Scatterplot with Groups

```
ggplot(data = df, aes(x = predictor, y = response)) +
  geom_point() + facet_wrap(~ categorical_variable, ncol = 2)
```

Vertical Bar Chart

```
ggplot(data = df, aes(x = col)) +
  geom_bar()
```

Grouped Bar Plot

```
ggplot(data = df, aes(x = col_1, fill = col_2)) +
  geom_bar(position = "dodge")
```

Explained

Plots a scatterplot with y-axis labels in comma format.

Creates scatterplots of response vs predictor, grouped by a categorical variable.

Creates scatterplots of response vs predictor, grouped by a categorical variable.

Creates a vertical bar chart to visualize counts of data.

Creates a grouped bar plot to compare frequency distributions of categorical variables.



Statistics & Probability

Syntax for

How to use

Mean

```
mean(x)
```

Median

```
median(x)
```

Weighted Mean

```
mean <- weighted.mean(x = distribution, w = weights)
```

Standard Deviation

```
sd(x)
```

Correlation

```
cor(x, y)
```

Linear Model

```
lm(y ~ x, data=df)
```

Types of Variables

```
# Example Variables: Age (Quantitative), Gender (Qualitative)
```

P-Value Decision Threshold

```
if (p_value < 0.05) { print('Reject null hypothesis') } else {
  print('Fail to reject null hypothesis') }
```

Explained

Calculates the mean of vector `x`.

Calculates the median of vector `x`.

Computes the weighted mean of a numerical vector using specific weights.

Calculates the standard deviation of `x`.

Calculates correlation between `x` and `y`.

Fits a linear regression model.

Classify variables as Quantitative (numerical) or Qualitative (categorical).

Decide on hypothesis rejection using a common p-value threshold of `0.05`.





Statistics & Probability

Syntax for	How to use	Explained	Syntax for	How to use	Explained
Chi-Squared Distribution	<pre>pchisq(3.84, df = 1)</pre>	Calculates the cumulative probability for a chi-squared distribution with specific degrees of freedom.	Simulate Coin Toss	<pre>set.seed(1) coin_toss <- function() { if (runif(1) <= 0.5) 'HEADS' else 'TAILS' }</pre>	Simulates a random coin toss using R's uniform random numbers.
Chi-Squared Test	<pre>pchisq(q = 10, df = 5)</pre>	Calculate cumulative probability for a chi-squared statistic of 10 with 5 degrees of freedom.	Addition Rule for Probability	$P(A \cup B) = P(A) + P(B) - P(A \cap B)$	Formula to calculate probabilities of unions of events, adjusting for overlap in non-exclusive cases.
Multi-category Chi-squared Test	<pre>data <- table(income\$sex, income\$high_income)</pre>	Performs a chi-squared test on the given contingency table.	Independent Events	$P(A \cap B) = P(A) * P(B)$	Probability of independent events occurs as product of individual probabilities.
Computing Mode in R	<pre>compute_mode <- function(vector) {counts_df <- tibble(vector) %>% group_by(vector) %>% summarise(frequency=n()) %>% arrange(desc(frequency)); counts_df\$vector[1]}</pre>	Defines a function to calculate the mode of a given vector using <code>dplyr</code> functions.	Product Rule in Experiments	<pre>total_outcomes <- a * b</pre>	Calculate the total outcomes for two independent experiments using the product rule.
Calculate Z-score	<pre>z_score <- function(value, vector) { (value - mean(vector)) / sd(vector) }</pre>	This calculates the Z-score for a value relative to a vector's distribution.	Uniform Distribution	<pre># Assuming all outcomes have equal chance outcomes <- c(1, 2, 3, 4, 5, 6) probabilities <- rep(1/6, 6) paste('Outcome:', outcomes, 'Probability:', probabilities)</pre>	Demonstrates a uniform distribution for a dice roll, where outcomes equally likely.
Chi-Squared Distribution	<pre>pchisq(3.84, df = 1)</pre>	Calculates the cumulative probability for a chi-squared distribution with specific degrees of freedom.			





Statistics & Probability

Syntax for	How to use	Explained
Conditional Probability Calculation	<pre>P_A_given_B <- P_A_and_B / P_B</pre>	Compute $P(A B)$ given the probability of A and B, and probability of B.
Conditional Probability	<pre>P_A_given_B <- length(intersect(A, B)) / length(B)</pre>	Compute $P(A B)$ using set cardinalities.
Conditional Probability Definition	<pre>P_A_given_B <- 1 - P_Ac_given_B</pre>	Conditional probabilities are interrelated; $P(A B)$ and its complement $P(A^c B)$ can be calculated mutually.
Independence	<pre>P_A_and_B <- P_A * P_B</pre>	Defines independent events: joint probability equals product of individual probabilities.



Programming

Syntax for	How to use	Explained
If Statement	<pre>if (x > 0) print("positive")</pre>	Executes code if condition is true.
For Loop	<pre>for (i in 1:3) print(i)</pre>	Iterates over a sequence.
While Loop	<pre>while (x < 5) x <- x + 1</pre>	Repeats code while the <code>x < 5</code> condition is true.
Syntax for functions	<pre>function_name <- function(input) { # Code to manipulate the input return(output) }</pre>	Defines a reusable function structure in R.
Define Function	<pre>f <- function(a, b) a + b</pre>	Defines a function with two arguments.
Apply Function	<pre>apply(m, 1, sum)</pre>	Applies a function over rows/columns of a matrix.
Exponentiation	<pre>3^5</pre>	Calculates 3 raised to the power of 5.
Creating Dates	<pre>ymd('20/04/21')</pre>	Converts a string into a Date object using 'year-month-day'.
Creating Dates from Strings	<pre>ymd("20/04/21")</pre>	Converts a string to a date object using the specified format.
Define Window Frame	<pre>ROWS BETWEEN 1 PRECEDING AND 1 FOLLOWING</pre>	Defines a window frame including one row before and after the current row for computations.





Machine Learning

Syntax for	How to use	Explained
Fitting a Linear Model	<pre>lm_fit <- lm(response ~ predictor, data = df)</pre>	Fit a linear regression model with a response and a predictor variable.
Visualize Residuals	<pre>library(ggplot2) ggplot(data.frame(residuals = lm_fit\$residuals), aes(x = residuals)) + geom_histogram()</pre>	Visualize the distribution of residuals to check the linear model's fit.
Hyperparameter Grid Search	<pre>knn_grid <- expand.grid(k = 1:20) knn_model <- train(tidy_price ~ accommodates + bathrooms + bedrooms, data = training_data, method = "knn", trControl = train_control, preProcess = c("center", "scale"), tuneGrid = knn_grid) plot(knn_model)</pre>	Performs grid search to optimize k for k-NN model and visualizes results.
Naive Bayes Algorithm	$P(\text{Spam} w_1, \dots, w_n) \propto P(\text{Spam}) * \prod_i P(w_i \text{Spam})$	Classifies messages as spam using conditional probabilities.



File I/O

Syntax for	How to use	Explained
Read CSV	<pre>read.csv("file.csv")</pre>	Reads a CSV file into a data frame.
Write CSV	<pre>write.csv(df, "file.csv")</pre>	Writes a data frame to a CSV file.
Read RDS	<pre>readRDS("file.rds")</pre>	Reads an RDS file into R.
Write RDS	<pre>saveRDS(df, "file.rds")</pre>	Saves an object as an RDS file.
List Files	<pre>list.files()</pre>	Lists files in the current directory.