

CheatSheet - Data Wrangling with Tidyverse



Commands	Syntax	Description	Example
install package	<code>install.packages("packagename")</code>	<code>install.packages</code> is used to install the packages from the R library.	<code>install.packages("tidyverse")</code>
load package	<code>library(packagename)</code>	<code>library()</code> Load the package from R library.	<code>library(tidyverse)</code>
download.file	<code>download.file(url, destfile, method, quiet = FALSE, mode = "w", cacheOK = TRUE, headers = NULL, ...)</code>	<code>download.file()</code> to download the file locally using the <code>download.file()</code> function. url naming the URL of a resource to be downloaded. destfile a character string with the name where the downloaded file is saved.	<code>download.file(url, destfile = "lax_to_jfk.tar.gz")</code>
untar	<code>untar()</code>	<code>untar()</code> is used to extract files from a tar archive is done with <code>untar</code> function from the <code>utils</code> package.	<code>untar("lax_to_jfk.tar.gz")</code>
read_csv	<code>read_csv(file)</code>	<code>read_csv()</code> reads the csv file using <code>readr</code> package.	<code>read_csv("lax_to_jfk/lax_to_jfk.csv")</code>
Missing Values and Formatting			
is.na	<code>is.na(x)</code>	<code>is.na(x)</code> returns a vector of TRUE or FALSE depending if the element in x is NA or not.	<code>is.na(c(1, na)) # FALSE TRUE</code>
anyNA	<code>anyNA(x, recursive = FALSE)</code>	<code>anyNA()</code> returns TRUE if x contains any NAs and FALSE otherwise.	<code>anyNA(c(1, na)) # TRUE</code>
sum	<code>sum(object)</code>	<code>sum()</code> is used to calculate sum.	<code>sum(is.na(carrierdelay))</code>
summarize	<code>summarize(X, by, FUN, ..., stat.name=deparse(substitute(X)), type=c('variables', 'matrix'), subset=TRUE, keepcolnames=FALSE)</code>	<code>summarize()</code> function reduces a data frame to a summary of just one vector or value. X a vector or matrix capable of being operated	<code>summarize(count = sum(is.na(carrierdelay)))</code>

		on by the function specified as the FUN argument	
		by one or more stratification variables. If a single variable, by may be a vector, otherwise it should be a list.	
		FUN a function of a single vector argument, used to create the statistical summaries for summarize. FUN may compute any number of statistics.	
map	map(.x, .f, ...)	map() functions transform their input by applying a function to each element and returning a vector the same length as the input.	map(sub_airline, ~sum(is.na(.)))
dim	dim(object)	dim returns the dimension of the matrix, array, or data frame.	dim(sub_airline)
drop_na	drop_na(object)	drop_na() drop rows containing missing values.	drop_na(carrierdelay)
		replace_na replace missing values.	
		data A data frame or vector.	
replace_na	replace_na(data, replace, ...)	replace If data is a data frame, a named list giving the value to replace NA with for each column. If data is a vector, a single value used for replacement.	replace_na(list(carrierdelay = 0, weatherdelay = 0, nasdelay = 0, securitydelay = 0, lateaircraftdelay = 0))
mean	mean(x, na.rm)	mean() calculate the arithmetic mean of the elements of the numeric vector passed to it as argument.	mean(drop_na_rows\$carrierdelay)
mutate, mutate_all, mutate_if	mutate(data, ...)	mutate function in R (mutate, mutate_all and mutate_at) is used to create new variable or column to the dataframe in R.	date_airline %>% select(year, month, day) %>% mutate_all(type.convert) %>% mutate_if(is.character, as.numeric)
Data Normalization			
Simple scaling	xnew=xold/xmax	Simple scaling divides each value by the maximum value in a feature. The	sub_airline\$arrdelay / max(sub_airline\$arrdelay)

Min-max	$x_{new} = (x_{old} - x_{max}) / (x_{max} - x_{min})$
---------	-------------------------------------------------------

Z-score	$x_{new} = (x_{old} - \mu) / \sigma$
---------	--------------------------------------

Binning Data

ggplot	<code>ggplot(df, aes(x, y, other aesthetics))</code>
--------	------------------------------------------------------

ntile	<code>ntile(data)</code>
-------	--------------------------

geom_histogram	<code>geom_histogram(*arguments)</code>
----------------	-----------------------------------------

Indicator variable

spread	<code>spread(data, key, value)</code>
--------	---------------------------------------

slice	<code>slice(num1 : num5)</code>
-------	----------------------------------

factor	<code>factor(x)</code>
--------	------------------------

new range is between 0 and 1.

Min-max subtracts the minimum value from the original and divides by the maximum minus the minimum. The minimum becomes 0 and the maximum becomes 1.

Standardization (Z-score)

subtracts the mean (μ) of the feature and divides by the standard deviation (σ).

```
(sub_airline$arrdelay -
min(sub_airline$arrdelay))
/(max(sub_airline$arrdelay) -
min(sub_airline$arrdelay))
```

```
(sub_airline$arrdelay -
mean(sub_airline$arrdelay)) /
sd(sub_airline$arrdelay)
```

ggplot is a plotting package that makes it simple to create complex plots from data in a data frame.

```
ggplot(data = sub_airline, mapping =
aes(x = arrdelay)) +
geom_histogram(bins = 100, color =
"white", fill = "red")
```

ntile() function is used to divide the data into N bins there by providing ntile rank.

```
sub_airline %>% mutate(quantile_rank
= ntile(sub_airline$arrdelay,4))
```

geom_histogram() function display the counts with bars.

```
geom_histogram(bins = 4, color =
"white", fill = "red")
```

spread a key-value pair across multiple columns
* data is your dataframe of interest.

* key is the column whose values will become variable names.

* value is the column where values will fill in under the new variables created from key.

slice() looks at the specified rows.

```
sub_airline %>%
spread(reporting_airline, arrdelay)
```

```
slice(1:5)
```

factor() function is used to encode a

vector as a factor, If argument ordered is TRUE, the factor levels are assumed to be ordered.

```
sub_airline %>%
mutate(reporting_airline =
factor(reporting_airline, labels =
c("aa", "as", "dl", "ua", "b6", "pa
(1)", "hp", "tw", "vx")))

```

Author(s)

[D.M. Naidu](#)

Changelog

Date	Version	Changed by	Change Description
2023-05-11	1.1	Eric Hao & Vladislav Boyko	Updated Page Frames
2020-08-11	1.0	D.M. Naidu	Initial Version