# Base R
## Cheat Sheet

## Getting Help

### Accessing the help files

**?mean**
Get help of a particular function.
**help.search('weighted mean')**
Search the help files for a word or phrase.
**help(package = 'dplyr')**
Find help for a package.

### More about an object

**str(iris)**
Get a summary of an object's structure.
**class(iris)**
Find the class an object belongs to.

## Using Packages

**install.packages('dplyr')**
Download and install a package from CRAN.

**library(dplyr)**
Load the package into the session, making all its functions available to use.

**dplyr::select**
Use a particular function from a package.

**data(iris)**
Load a built-in dataset into the environment.

## Working Directory

**getwd()**
Find the current working directory (where inputs are found and outputs are sent).

**setwd('C://file/path')**
Change the current working directory.

**Use projects in RStudio to set the working directory to the folder you are working in.**

## Vectors

### Creating Vectors

| | | |
|---|---|---|
| `c(2, 4, 6)` | 2 4 6 | Join elements into a vector |
| `2:6` | 2 3 4 5 6 | An integer sequence |
| `seq(2, 3, by=0.5)` | 2.0 2.5 3.0 | A complex sequence |
| `rep(1:2, times=3)` | 1 2 1 2 1 2 | Repeat a vector |
| `rep(1:2, each=3)` | 1 1 1 2 2 2 | Repeat elements of a vector |

### Vector Functions

**sort(x)**
Return x sorted.
**table(x)**
See counts of values.
**rev(x)**
Return x reversed.
**unique(x)**
See unique values.

### Selecting Vector Elements

#### By Position

| | |
|---|---|
| **x[4]** | The fourth element. |
| **x[-4]** | All but the fourth. |
| **x[2:4]** | Elements two to four. |
| **x[-(2:4)]** | All elements except two to four. |
| **x[c(1, 5)]** | Elements one and five. |

#### By Value

| | |
|---|---|
| **x[x == 10]** | Elements which are equal to 10. |
| **x[x < 0]** | All elements less than zero. |
| **x[x %in% c(1, 2, 5)]** | Elements in the set 1, 2, 5. |

#### Named Vectors

| | |
|---|---|
| **x['apple']** | Element with name 'apple'. |

## Programming

### For Loop

```
for (variable in sequence){
    Do something
}
```

#### Example

```
for (i in 1:4){
    j <- i + 10
    print(j)
}
```

### While Loop

```
while (condition){
    Do something
}
```

#### Example

```
while (i < 5){
    print(i)
    i <- i + 1
}
```

### If Statements

```
if (condition){
    Do something
} else {
    Do something different
}
```

#### Example

```
if (i > 3){
    print('Yes')
} else {
    print('No')
}
```

### Functions

```
function_name <- function(var){
    Do something
    return(new_variable)
}
```

#### Example

```
square <- function(x){
    squared <- x*x
    return(squared)
}
```

### Reading and Writing Data

Also see the **readr** package.

| Input | Ouput | Description |
|---|---|---|
| `df <- read.table('file.txt')` | `write.table(df, 'file.txt')` | Read and write a delimited text file. |
| `df <- read.csv('file.csv')` | `write.csv(df, 'file.csv')` | Read and write a comma separated value file. This is a special case of read.table/write.table. |
| `load('file.RData')` | `save(df, file = 'file.Rdata')` | Read and write an R data file, a file type special for R. |

| Conditions | | | | | | | |
|---|---|---|---|---|---|---|---|
| `a == b` | Are equal | `a > b` | Greater than | `a >= b` | Greater than or equal to | `is.na(a)` | Is missing |
| `a != b` | Not equal | `a < b` | Less than | `a <= b` | Less than or equal to | `is.null(a)` | Is null |

# Types

Converting between common data types in R. Can always go from a higher value in the table to a lower value.

| | | |
|---|---|---|
| as.logical | TRUE, FALSE, TRUE | Boolean values (TRUE or FALSE). |
| as.numeric | 1, 0, 1 | Integers or floating point numbers. |
| as.character | '1', '0', '1' | Character strings. Generally preferred to factors. |
| as.factor | '1', '0', '1', levels: '1', '0' | Character strings with preset levels. Needed for some statistical models. |

# Maths Functions

| | | | | |
|---|---|---|---|---|
| log(x) | Natural log. | sum(x) | Sum. |
| exp(x) | Exponential. | mean(x) | Mean. |
| max(x) | Largest element. | median(x) | Median. |
| min(x) | Smallest element. | quantile(x) | Percentage quantiles. |
| round(x, n) | Round to n decimal places. | rank(x) | Rank of elements. |
| signif(x, n) | Round to n significant figures. | var(x) | The variance. |
| cor(x, y) | Correlation. | sd(x) | The standard deviation. |

# Variable Assignment

```
> a <- 'apple'
> a
[1] 'apple'
```

## The Environment

| | |
|---|---|
| ls() | List all variables in the environment. |
| rm(x) | Remove x from the environment. |
| rm(list = ls()) | Remove all variables from the environment. |

**You can use the environment panel in RStudio to browse variables in your environment.**

# Matrices

```
m <- matrix(x, nrow = 3, ncol = 3)
```
Create a matrix from x.

m[2, ] - Select a row

m[ , 1] - Select a column

m[2, 3] - Select an element

t(m)
Transpose

m %*% n
Matrix Multiplication

solve(m, n)
Find x in: m * x = n

# Lists

```
l <- list(x = 1:5, y = c('a', 'b'))
```
A list is a collection of elements which can be of different types.

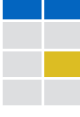| l[[2]] | l[1] | l$x | l['y'] |
|---|---|---|---|
| Second element of l. | New list with only the first element. | Element named x. | New list with only element named y. |

# Data Frames

Also see the **dplyr** package.

```
df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))
```
A special case of a list where all elements are the same length.

| x | y |
|---|---|
| 1 | a |
| 2 | b |
| 3 | c |

## List subsetting

df$x     df[[2]]

*Understanding a data frame*

View(df) — See the full data frame.

head(df) — See the first 6 rows.

## Matrix subsetting
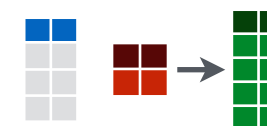
df[ , 2]

df[2, ]

df[2, 2]

nrow(df)
Number of rows.

ncol(df)
Number of columns.

dim(df)
Number of columns and rows.

cbind - Bind columns.

rbind - Bind rows.

# Strings

Also see the **stringr** package.

| | |
|---|---|
| paste(x, y, sep = ' ') | Join multiple vectors together. |
| paste(x, collapse = ' ') | Join elements of a vector together. |
| grep(pattern, x) | Find regular expression matches in x. |
| gsub(pattern, replace, x) | Replace matches in x with a string. |
| toupper(x) | Convert to uppercase. |
| tolower(x) | Convert to lowercase. |
| nchar(x) | Number of characters in a string. |

# Factors

factor(x)
Turn a vector into a factor. Can set the levels of the factor and the order.

cut(x, breaks = 4)
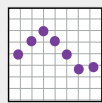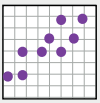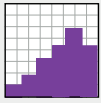Turn a numeric vector into a factor by 'cutting' into sections.

# Statistics

lm(y ~ x, data=df)
Linear model.

glm(y ~ x, data=df)
Generalised linear model.

summary
Get more detailed information out a model.

t.test(x, y)
Perform a t-test for difference between means.

pairwise.t.test
Perform a t-test for paired data.

prop.test
Test for a difference between proportions.

aov
Analysis of variance.

# Distributions

| | Random Variates | Density Function | Cumulative Distribution | Quantile |
|---|---|---|---|---|
| Normal | rnorm | dnorm | pnorm | qnorm |
| Poisson | rpois | dpois | ppois | qpois |
| Binomial | rbinom | dbinom | pbinom | qbinom |
| Uniform | runif | dunif | punif | qunif |

# Plotting

Also see the **ggplot2** package.

plot(x)
Values of x in order.

plot(x, y)
Values of x against y.

hist(x)
Histogram of x.

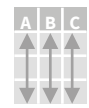# Dates

See the **lubridate** package.

# Data transformation with dplyr : : **CHEATSHEET**

**dplyr** functions work with pipes and expect **tidy data**. In tidy data:

Each **variable** is in its own **column**

Each **observation**, or **case**, is in its own **row**

**pipes**

x |> f(y)
becomes f(x, y)

## Summarize Cases

Apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

**summary function**

**summarize(**.data, …**)**
Compute table of summaries.
mtcars |> summarize(avg = mean(mpg))

**count(**.data, …, wt = NULL, sort = FALSE, name = NULL**)** Count number of rows in each group defined by the variables in … Also **tally()**, **add_count()**, **add_tally()**.
mtcars |> count(cyl)

## Group Cases

Use **group_by(**.data, …, .add = FALSE, .drop = TRUE**)** to create a "grouped" copy of a table grouped by columns in … dplyr functions will manipulate each "group" separately and combine the results.

mtcars |>
  group_by(cyl) |>
  summarize(avg = mean(mpg))

Use **rowwise(**.data, …**)** to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidyr cheat sheet for list-column workflow.

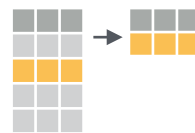starwars |>
  rowwise() |>
  mutate(film_count = length(films))

**ungroup(**x, …**)** Returns ungrouped copy of table.
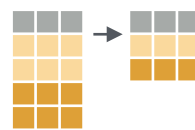g_mtcars <- mtcars |> group_by(cyl)
ungroup(g_mtcars)

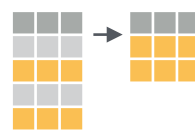## Manipulate Cases

### EXTRACT CASES
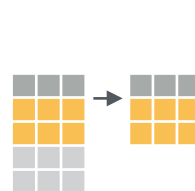
Row functions return a subset of rows as a new table.

**filter(**.data, …, .preserve = FALSE**)** Extract rows that meet logical criteria.
mtcars |> filter(mpg > 20)

**distinct(**.data, …, .keep_all = FALSE**)** Remove rows with duplicate values.
mtcars |> distinct(gear)

**slice(**.data, …, .preserve = FALSE**)** Select rows by position.
mtcars |> slice(10:15)

**slice_sample(**.data, …, n, prop, weight_by = NULL, replace = FALSE**)** Randomly select rows. Use n to select a number of rows and prop to select a fraction of rows.
mtcars |> slice_sample(n = 5, replace = TRUE)

**slice_min(**.data, order_by, …, n, prop, with_ties = TRUE**)** and **slice_max()** Select rows with the lowest and highest values.
mtcars |> slice_min(mpg, prop = 0.25)

**slice_head(**.data, …, n, prop**)** and **slice_tail()** Select the first or last rows.
mtcars |> slice_head(n = 5)

**Logical and boolean operators to use with filter()**

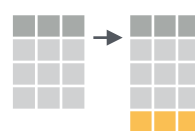| == | < | <= | is.na() | %in% | \| | xor() |
| != | > | >= | !is.na() | ! | & | |

See **?base::Logic** and **?Comparison** for help.

### ARRANGE CASES

**arrange(**.data, …, .by_group = FALSE**)** Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low.
mtcars |> arrange(mpg)
mtcars |> arrange(desc(mpg))

### ADD CASES

**add_row(**.data, …, .before = NULL, .after = NULL**)** Add one or more rows to a table.
cars |> add_row(speed = 1, dist = 1)

## Manipulate Variables

### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

**pull(**.data, var = -1, name = NULL, …**)** Extract column values as a vector, by name or index.
mtcars |> pull(wt)

**select(**.data, …**)** Extract columns as a table.
mtcars |> select(mpg, wt)

**relocate(**.data, …, .before = NULL, .after = NULL**)** Move columns to new position.
mtcars |> relocate(mpg, cyl, .after = last_col())

**Use these helpers with select() and across()**
e.g. mtcars |> select(mpg:cyl)

| **contains(**match**)** | **num_range(**prefix, range**)** | **:**, e.g., mpg:cyl |
| **ends_with(**match**)** | **all_of(**x**)**/**any_of(**x, …, vars**)** | **!**, e.g., !gear |
| **starts_with(**match**)** | **matches(**match**)** | **everything()** |

### MANIPULATE MULTIPLE VARIABLES AT ONCE

df <- tibble(x_1 = c(1, 2), x_2 = c(3, 4), y = c(4, 5))

**across(**.cols, .funs, …, .names = NULL**)** Summarize or mutate multiple columns in the same way.
df |> summarize(across(everything(), mean))

**c_across(**.cols**)** Compute across columns in row-wise data.
df |>
  rowwise() |>
  mutate(x_total = sum(c_across(1:2)))

### MAKE NEW VARIABLES

Apply **vectorized functions** to columns. Vectorized functions take vectors as input and return vectors of the same length as output (see back).

**vectorized function**

**mutate(**.data, …, .keep = "all", .before = NULL, .after = NULL**)** Compute new column(s). Also **add_column()**.
mtcars |> mutate(gpm = 1 / mpg)
mtcars |> mutate(gpm = 1 / mpg, .keep = "none")

**rename(**.data, …**)** Rename columns. Use **rename_with()** to rename with a function.
mtcars |> rename(miles_per_gallon = mpg)

posit®

# Vectorized Functions

## TO USE WITH MUTATE ()

**mutate()** applies vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

> **vectorized function**

### OFFSET
dplyr::**lag()** - offset elements by 1
dplyr::**lead()** - offset elements by -1

### CUMULATIVE AGGREGATE
dplyr::**cumall()** - cumulative all()
dplyr::**cumany()** - cumulative any()
**cummax()** - cumulative max()
dplyr::**cummean()** - cumulative mean()
**cummin()** - cumulative min()
**cumprod()** - cumulative prod()
**cumsum()** - cumulative sum()

### RANKING
dplyr::**cume_dist()** - proportion of all values <=
dplyr::**dense_rank()** - rank w ties = min, no gaps
dplyr::**min_rank()** - rank with ties = min
dplyr::**ntile()** - bins into n bins
dplyr::**percent_rank()** - min_rank scaled to [0,1]
dplyr::**row_number()** - rank with ties = "first"

### MATH
**+, -, \*, /, ^, %/%, %%** - arithmetic ops
**log(), log2(), log10()** - logs
**<, <=, >, >=, !=, ==** - logical comparisons
dplyr::**between()** - x >= left & x <= right
dplyr::**near()** - safe == for floating point numbers

### MISCELLANEOUS
dplyr::**case_when()** - multi-case if_else()
```
starwars |>
    mutate(type = case_when(
        height > 200 | mass > 200 ~ "large",
        species == "Droid"        ~ "robot",
        TRUE                      ~ "other")
    )
```
dplyr::**coalesce()** - first non-NA values by element across a set of vectors
dplyr::**if_else()** - element-wise if() + else()
dplyr::**na_if()** - replace specific values with NA
**pmax()** - element-wise max()
**pmin()** - element-wise min()

# Summary Functions

## TO USE WITH SUMMARIZE ()

**summarize()** applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

> **summary function**

### COUNT
dplyr::**n()** - number of values/rows
dplyr::**n_distinct()** - # of uniques
**sum(!is.na())** - # of non-NAs

### POSITION
**mean()** - mean, also **mean(!is.na())**
**median()** - median

### LOGICAL
**mean()** - proportion of TRUEs
**sum()** - # of TRUEs

### ORDER
dplyr::**first()** - first value
dplyr::**last()** - last value
dplyr::**nth()** - value in nth location of vector

### RANK
**quantile()** - nth quantile
**min()** - minimum value
**max()** - maximum value

### SPREAD
**IQR()** - Inter-Quartile Range
**mad()** - median absolute deviation
**sd()** - standard deviation
**var()** - variance

# Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

tibble::**rownames_to_column()**
Move row names into col.
```
a <- mtcars |>
    rownames_to_column(var = "C")
```

tibble::**column_to_rownames()**
Move col into row names.
```
a |> column_to_rownames(var = "C")
```

Also tibble::**has_rownames()** and tibble::**remove_rownames()**.

# Combine Tables

## COMBINE VARIABLES

**bind_cols**(…, .name_repair**)** Returns tables placed side by side as a single table. Column lengths must be equal. Columns will NOT be matched by id (to do that look at Relational Data below), so be sure to check that both tables are ordered the way you want before binding.

## RELATIONAL DATA

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

**left_join**(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join matching values from y to x.

**right_join**(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join matching values from x to y.

**inner_join**(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join data. Retain only rows with matches.

**full_join**(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join data. Retain all values, all rows.

## COLUMN MATCHING FOR JOINS

Use **by = c("col1", "col2", …)** to specify one or more common columns to match on.
left_join(x, y, by = "A")

Use a named vector, **by = c("col1" = "col2")**, to match on columns that have different names in each table.
left_join(x, y, by = c("C" = "D"))

Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = c("C" = "D"),
suffix = c("1", "2"))

## COMBINE CASES

**bind_rows**(…, .id = NULL**)** Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured).

Use a "**Filtering Join**" to filter one table against the rows of another.

**semi_join**(x, y, by = NULL, copy = FALSE, …, na_matches = "na"**)** Return rows of x that have a match in y. Use to see what will be included in a join.

**anti_join**(x, y, by = NULL, copy = FALSE, …, na_matches = "na"**)** Return rows of x that do not have a match in y. Use to see what will not be included in a join.

Use a "**Nest Join**" to inner join one table to another into a nested data frame.

**nest_join**(x, y, by = NULL, copy = FALSE, keep = FALSE, name = NULL, …**)** Join data, nesting matches from y in a single new data frame column.

## SET OPERATIONS

**intersect(x, y, …)**
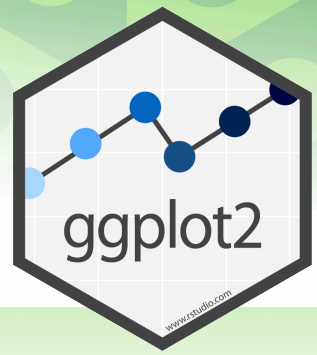Rows that appear in both x and y.

**setdiff(x, y, …)**
Rows that appear in x but not y.

**union(x, y, …)**
Rows that appear in x or y, duplicates removed). **union_all()** retains duplicates.

Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

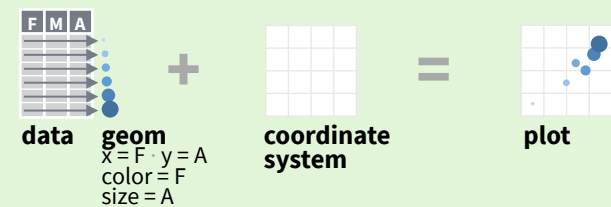# Data visualization with ggplot2 : : **CHEATSHEET**

ggplot2

## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data** set, a **coordinate system**, and **geoms**—visual marks that represent data points.



data | geom
x = F · y = A
+ coordinate system = plot

To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.

data | geom
x = F · y = A
color = F
size = A
+ coordinate system = plot

Complete the template below to build a graph.

**ggplot (data = `<DATA>`) +**  `required`
**`<GEOM_FUNCTION>`(mapping = aes(`<MAPPINGS>`),**
  stat = `<STAT>`, position = `<POSITION>`) +   `Not required, sensible defaults supplied`
  `<COORDINATE_FUNCTION>` +
  `<FACET_FUNCTION>` +
  `<SCALE_FUNCTION>` +
  `<THEME_FUNCTION>`

**ggplot**(data = mpg, **aes**(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

**last_plot()** Returns the last plot.

**ggsave**("plot.png", width = 5, height = 5) Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

## Aes   Common aesthetic values.

**color** and **fill** - string ("red", "#RRGGBB")

**linetype** - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

**size** - integer (in mm for size of points and text)

**linewidth** - integer (in mm for widths of lines)

**shape** - integer/shape name or a single character ("a")

0 1 2 3 4 5 6 7 8 9 10 11 12
□ ○ △ + × ◇ ▽ ⊠ * ⊕ ⊞ ⊠ ⊞
13 14 15 16 17 18 19 20 21 22 23 24 25
⊠ ◻ ■ ● ▲ ◆ ● ● ● ◆ ◇ △ ▽

## Geoms  Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))

**a + geom_blank()** and **a + expand_limits()** Ensure limits include values across all plots.

**b + geom_curve**(aes(yend = lat + 1, xend = long + 1), curvature = 1) - x, xend, y, yend, alpha, angle, color, curvature, linetype, size

**a + geom_path**(lineend = "butt", linejoin = "round", linemitre = 1) x, y, alpha, color, group, linetype, size

**a + geom_polygon**(aes(alpha = 50)) - x, y, alpha, color, fill, group, subgroup, linetype, size

**b + geom_rect**(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1)) - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size

**a + geom_ribbon**(aes(ymin = unemploy - 900, ymax = unemploy + 900)) - x, ymax, ymin, alpha, color, fill, group, linetype, size

### LINE SEGMENTS
common aesthetics: x, y, alpha, color, linetype, size

**b + geom_abline**(aes(intercept = 0, slope = 1))
**b + geom_hline**(aes(yintercept = lat))
**b + geom_vline**(aes(xintercept = long))
**b + geom_segment**(aes(yend = lat + 1, xend = long + 1))
**b + geom_spoke**(aes(angle = 1:1155, radius = 1))

### ONE VARIABLE    continuous
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

**c + geom_area**(stat = "bin") x, y, alpha, color, fill, linetype, size

**c + geom_density**(kernel = "gaussian") x, y, alpha, color, fill, group, linetype, size, weight

**c + geom_dotplot()** x, y, alpha, color, fill

**c + geom_freqpoly()** x, y, alpha, color, group, linetype, size

**c + geom_histogram**(binwidth = 5) x, y, alpha, color, fill, linetype, size, weight

**c2 + geom_qq**(aes(sample = hwy)) x, y, alpha, color, fill, linetype, size, weight

### discrete
d <- ggplot(mpg, aes(fl))

**d + geom_bar()** x, alpha, color, fill, linetype, size, weight

### TWO VARIABLES
#### both continuous
e <- ggplot(mpg, aes(cty, hwy))

**e + geom_label**(aes(label = cty), nudge_x = 1, nudge_y = 1) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

**e + geom_point()** x, y, alpha, color, fill, shape, size, stroke

**e + geom_quantile()** x, y, alpha, color, group, linetype, size, weight

**e + geom_rug**(sides = "bl") x, y, alpha, color, linetype, size

**e + geom_smooth**(method = lm) x, y, alpha, color, fill, group, linetype, size, weight

**e + geom_text**(aes(label = cty), nudge_x = 1, nudge_y = 1) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

#### one discrete, one continuous
f <- ggplot(mpg, aes(class, hwy))

**f + geom_col()** x, y, alpha, color, fill, group, linetype, size

**f + geom_boxplot()** x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

**f + geom_dotplot**(binaxis = "y", stackdir = "center") x, y, alpha, color, fill, group

**f + geom_violin**(scale = "area") x, y, alpha, color, fill, group, linetype, size, weight

#### both discrete
g <- ggplot(diamonds, aes(cut, color))

**g + geom_count()** x, y, alpha, color, fill, shape, size, stroke

**e + geom_jitter**(height = 2, width = 2) x, y, alpha, color, fill, shape, size

### continuous bivariate distribution
h <- ggplot(diamonds, aes(carat, price))

**h + geom_bin2d**(binwidth = c(0.25, 500)) x, y, alpha, color, fill, linetype, size, weight

**h + geom_density_2d()** x, y, alpha, color, group, linetype, size

**h + geom_hex()** x, y, alpha, color, fill, size

### continuous function
i <- ggplot(economics, aes(date, unemploy))

**i + geom_area()** x, y, alpha, color, fill, linetype, size

**i + geom_line()** x, y, alpha, color, group, linetype, size

**i + geom_step**(direction = "hv") x, y, alpha, color, group, linetype, size

### visualizing error
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))

**j + geom_crossbar**(fatten = 2) - x, y, ymax, ymin, alpha, color, fill, group, linetype, size

**j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width Also **geom_errorbarh()**.

**j + geom_linerange()** x, ymin, ymax, alpha, color, group, linetype, size

**j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

### maps
Draw the appropriate geometric object depending on the simple features present in the data. aes() arguments: map_id, alpha, color, fill, linetype, linewidth.

nc <- **sf::st_read**(system.file("shape/nc.shp", package = "sf"))

ggplot(nc) +
  **geom_sf**(aes(fill = AREA))

### THREE VARIABLES
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)); l <- ggplot(seals, aes(long, lat))

**l + geom_contour**(aes(z = z)) x, y, z, alpha, color, group, linetype, size, weight

**l + geom_contour_filled**(aes(fill = z)) x, y, alpha, color, fill, group, linetype, size, subgroup
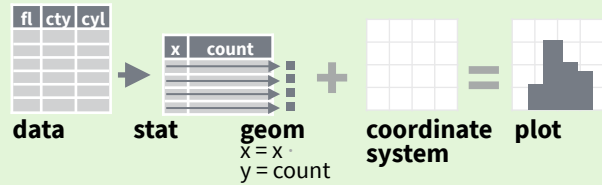
**l + geom_raster**(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE) x, y, alpha, fill

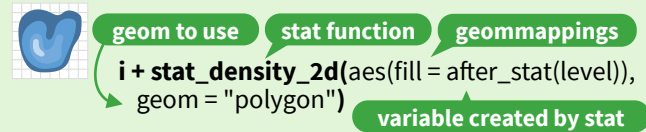**l + geom_tile**(aes(fill = z)) x, y, alpha, color, fill, linetype, size, width

**posit**®

# Stats An alternative way to build a layer.

A stat builds new variables to plot (e.g., count, prop).



| data | stat | geom | coordinate system | plot |

x = x ·
y = count

Visualize a stat by changing the default stat of a geom function, **geom_bar(stat="count")** or by using a stat function, **stat_count(geom="bar")**, which calls a default geom to make a layer (equivalent to a geom function).
Use **after_stat(name)** syntax to map the stat variable **name** to an aesthetic.

**geom to use** **stat function** **geommappings**
**i + stat_density_2d(**aes(fill = after_stat(level)),
geom = "polygon"**)** **variable created by stat**

**c + stat_bin(**binwidth = 1, boundary = 10**)**
**x, y** │ count, ncount, density, ndensity

**c + stat_count(**width = 1**) x, y** │ count, prop

**c + stat_density(**adjust = 1, kernel = "gaussian"**)**
**x, y** │ count, density, scaled

**e + stat_bin_2d(**bins = 30, drop = T**)**
**x, y, fill** │ count, density

**e + stat_bin_hex(**bins = 30**) x, y, fill** │ count, density

**e + stat_density_2d(**contour = TRUE, n = 100**)**
**x, y, color, size** │ level

**e + stat_ellipse(**level = 0.95, segments = 51, type = "t"**)**

**l + stat_contour(**aes(z = z)**) x, y, z, order** │ level

**l + stat_summary_hex(**aes(z = z), bins = 30, fun = max**)**
**x, y, z, fill** │ value

**l + stat_summary_2d(**aes(z = z), bins = 30, fun = mean**)**
**x, y, z, fill** │ value

**f + stat_boxplot(**coef = 1.5**)**
**x, y** │ lower, middle, upper, width , ymin, ymax

**f + stat_ydensity(**kernel = "gaussian", scale = "area"**) x, y** │
density, scaled, count, n, violinwidth, width

**e + stat_ecdf(**n = 40**) x, y** │ x, y

**e + stat_quantile(**quantiles = c(0.1, 0.9),
formula = y ~ log(x), method = "rq"**) x, y** │ quantile

**e + stat_smooth(**method = "lm", formula = y ~ x, se = T,
level = 0.95**) x, y** │ se, x, y, ymin, ymax

**ggplot() + xlim(**-5, 5**) + stat_function(**fun = dnorm,
n = 20, geom = "point"**) x** │ x, y

**ggplot() + stat_qq(**aes(sample = 1:100)**)**
**x, y, sample** │ sample, theoretical

**e + stat_sum() x, y, size** │ n, prop

**e + stat_summary(**fun.data = "mean_cl_boot"**)**

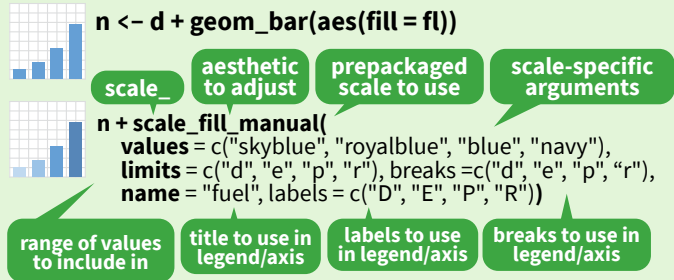**h + stat_summary_bin(**fun = "mean", geom = "bar"**)**

**e + stat_identity()**

**e + stat_unique()**

# Scales Override defaults with **scales** package.

**Scales** map data values to the visual values of an aesthetic. To change a mapping, add a new scale.

**n <- d + geom_bar(aes(fill = fl))**

**scale_** **aesthetic to adjust** **prepackaged scale to use** **scale-specific arguments**

**n + scale_fill_manual(**
**values** = c("skyblue", "royalblue", "blue", "navy"),
**limits** = c("d", "e", "p", "r"), breaks =c("d", "e", "p", "r"),
**name** = "fuel", labels = c("D", "E", "P", "R")**)**

**range of values to include** **title to use in legend/axis** **labels to use in legend/axis** **breaks to use in legend/axis**

## GENERAL PURPOSE SCALES

Use with most aesthetics

**scale_*_continuous()** - Map cont' values to visual ones.
**scale_*_discrete()** - Map discrete values to visual ones.
**scale_*_binned()** - Map continuous values to discrete bins.
**scale_*_identity()** - Use data values as visual ones.
**scale_*_manual(**values = c()**)** - Map discrete values to manually chosen visual ones.
**scale_*_date(**date_labels = "%m/%d"),
date_breaks = "2 weeks"**)** - Treat data values as dates.
**scale_*_datetime()** - Treat data values as date times.
Same as scale_*_date(). See ?strptime for label formats.

## X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

**scale_x_log10()** - Plot x on log10 scale.
**scale_x_reverse()** - Reverse the direction of the x axis.
**scale_x_sqrt()** - Plot x on square root scale.

## COLOR AND FILL SCALES (DISCRETE)

**n + scale_fill_brewer(**palette = "Blues"**)**
For palette choices:
RColorBrewer::display.brewer.all()

**n + scale_fill_grey(**start = 0.2,
end = 0.8, na.value = "red"**)**

## COLOR AND FILL SCALES (CONTINUOUS)

**o <- c + geom_dotplot(**aes(fill = x)**)**

**o + scale_fill_distiller(**palette = "Blues"**)**

**o + scale_fill_gradient(**low="red", high="yellow"**)**

**o + scale_fill_gradient2(**low = "red", high = "blue",
mid = "white", midpoint = 25**)**

**o + scale_fill_gradientn(**colors = topo.colors(6)**)**
Also: rainbow(), heat.colors(), terrain.colors(),
cm.colors(), RColorBrewer::brewer.pal()

## SHAPE AND SIZE SCALES

p <- e + geom_point(aes(shape = fl, size = cyl))

**p + scale_shape() + scale_size()**
**p + scale_shape_manual(**values = c(3:7)**)**

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
□ ○ △ + × ◇ ▽ ✳ ❋ ⊕ ⊞ ⊠ ⊗ □ △ ○ ○ ● ● ◆ ▲ ▼

**p + scale_radius(**range = c(1,6)**)**
**p + scale_size_area(**max_size = 6**)**

# Coordinate Systems

r <- d + geom_bar()

**r + coord_cartesian(**xlim = c(0, 5)**)** - xlim, ylim
The default cartesian coordinate system.

**r + coord_fixed(**ratio = 1/2**)**
ratio, xlim, ylim - Cartesian coordinates with fixed aspect ratio between x and y units.

**r + coord_flip()**
Flip cartesian coordinates by switching x and y aesthetic mappings.

**r + coord_polar(**theta = "x", direction=1**)**
theta, start, direction - Polar coordinates.

**r + coord_trans(**y = "sqrt"**)** - x, y, xlim, ylim
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

**π + coord_sf()** - xlim, ylim, crs. Ensures all layers use a common Coordinate Reference System.

# Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

s <- ggplot(mpg, aes(fl, fill = drv))

**s + geom_bar(position = "dodge")**
Arrange elements side by side.

**s + geom_bar(position = "fill")**
Stack elements on top of one another, normalize height.

**e + geom_point(position = "jitter")**
Add random noise to X and Y position of each element to avoid overplotting.

**e + geom_label(position = "nudge")**
Nudge labels away from points.

**s + geom_bar(position = "stack")**
Stack elements on top of one another.

Each position adjustment can be recast as a function with manual **width** and **height** arguments:
s + geom_bar(position = position_dodge(width = 1))

# Themes

**r + theme_bw()**
White background with grid lines.

**r + theme_classic()**

**r + theme_light()**

**r + theme_gray()**
Grey background (default theme).

**r + theme_linedraw()**

**r + theme_minimal()**
Minimal theme.

**r + theme_dark()**
Dark for contrast.

**r + theme_void()**
Empty theme.

**r + theme()** Customize aspects of the theme such as axis, legend, panel, and facet properties.
r + labs(title = "Title") + theme(plot.title.position = "plot")
r + theme(panel.background = element_rect(fill = "blue"))

# Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

t <- ggplot(mpg, aes(cty, hwy)) + geom_point()

**t + facet_grid(. ~ fl)**
Facet into columns based on fl.

**t + facet_grid(year ~ .)**
Facet into rows based on year.

**t + facet_grid(year ~ fl)**
Facet into both rows and columns.

**t + facet_wrap(~ fl)**
Wrap facets into a rectangular layout.

Set **scales** to let axis limits vary across facets.

**t + facet_grid(drv ~ fl, scales = "free")**
x and y axis limits adjust to individual facets:
**"free_x"** - x axis limits adjust
**"free_y"** - y axis limits adjust

Set **labeller** to adjust facet label:

**t + facet_grid(. ~ fl, labeller = label_both)**

| fl: c | fl: d | fl: e | fl: p | fl: r |

**t + facet_grid(fl ~ ., labeller = label_bquote(**alpha ^ .(fl)**))**

| $\alpha^c$ | $\alpha^d$ | $\alpha^e$ | $\alpha^p$ | $\alpha^r$ |

# Labels and Legends

Use **labs()** to label the elements of your plot.

**t + labs(x** = "New x axis label", **y** = "New y axis label",
**title** ="Add a title above the plot",
**subtitle** = "Add a subtitle below title",
**caption** = "Add a caption below plot",
**alt** = "Add alt text to the plot",
**<AES>** = "New **<AES>** legend title"**)**

**t + annotate(**geom = "text", x = 8, y = 9, label = "A"**)**
Places a geom with manually selected aesthetics.

**p + guides(**x = guide_axis(n.dodge = 2)**)** Avoid crowded or overlapping labels with guide_axis(n.dodge or angle).

**n + guides(**fill = "none"**)** Set legend type for each aesthetic: colorbar, legend, or none (no legend).

**n + theme(**legend.position = "bottom"**)**
Place legend at "bottom", "top", "left", or "right".

**n + scale_fill_discrete(**name = "Title",
labels = c("A", "B", "C", "D", "E")**)**
Set legend title and labels with a scale function.

# Zooming

**Without clipping** (preferred):
**t + coord_cartesian(**xlim = c(0, 100), ylim = c(10, 20)**)**

**With clipping** (removes unseen data points):
**t + xlim(**0, 100**) + ylim(**10, 20**)**

**t + scale_x_continuous(**limits = c(0, 100)**) +**
**scale_y_continuous(**limits = c(0, 100)**)**

# R GRAPHICAL PARAMETERS CHEATSHEET

## Box feature

**bty** : kind of box

*o=complete / ?=top & right / n=no box / c=top & left & bottom / l=bottom & left*

## Title

**main** : name of the title
**cex.main** : size *cex.main=2*
**col.main** : color *col.main="red"*
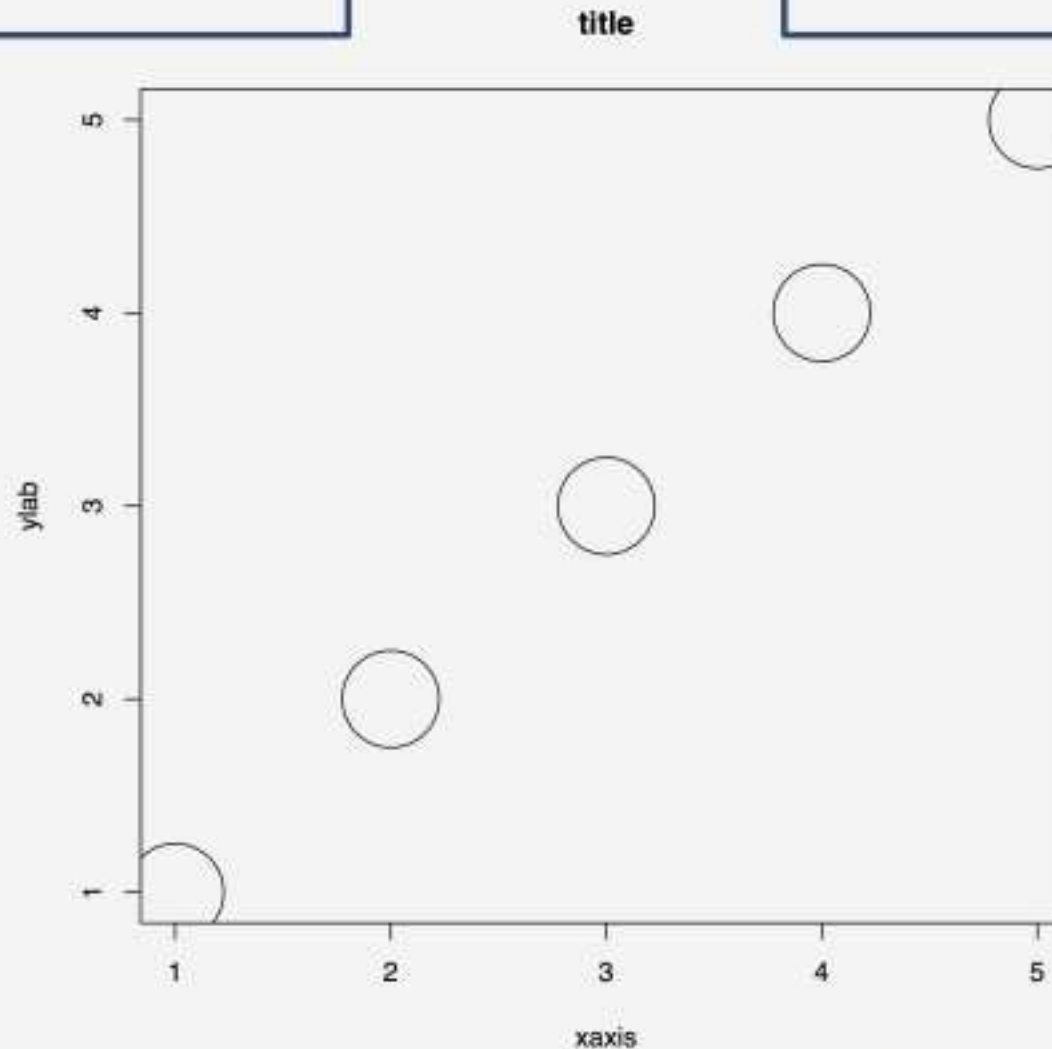**font.main** : font *font.main=3*

title

## Margins

See **Graph #135**
*mar, oma, omd, omi*

## Symbol styles

See **Graph #6**
*pch, lwd, pty,*
*col, cex, type...*

## General

**tck** : add a grid on a % of the area *tck=1*
**bg** : color of the background *par(bg="red")*
**font** : font of the text *(normal, bold, italic..)*
**lheight**: size between lines of titles
**srt**: text rotation in degree

ylab

xaxis

## X-axis name

**xlab** : name of the axis
**cex.lab** : size *cex.lab=2*
**col.lab** : color *col.lab="red"*
**sub** : to add a subtitle

## X-axis features

**lab** : number of graduation *lab=c(12,2,0)*
**xaxp** : to add c graduation from a to b: *xaxp=c(a,b,c)*
**log** : for logarithmic scale: *log="x"*
**xaxt** : to remove x axis: *xaxt="n"*
**fg** : color of axis, ticks and grid: *fg="red"*
**cex.axis** : size of tick labels
**col.axis** : color of tick labels
**xlim** : limits of the axis *xlim=c(0,10)*
**las**: orientation of tick labels
*0=parralel to the axis, 1=horizontal...*

| | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
| 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 |
| 43 | 44 | 45 | 46 | 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 |
| 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 | 80 | 81 | 82 | 83 | 84 |
| 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 | 101 | 102 | 103 | 104 | 105 |
| 106 | 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 | 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 |
| 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 | 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 | 147 |
| 148 | 149 | 150 | 151 | 152 | 153 | 154 | 155 | 156 | 157 | 158 | 159 | 160 | 161 | 162 | 163 | 164 | 165 | 166 | 167 | 168 |
| 169 | 170 | 171 | 172 | 173 | 174 | 175 | 176 | 177 | 178 | 179 | 180 | 181 | 182 | 183 | 184 | 185 | 186 | 187 | 188 | 189 |
| 190 | 191 | 192 | 193 | 194 | 195 | 196 | 197 | 198 | 199 | 200 | 201 | 202 | 203 | 204 | 205 | 206 | 207 | 208 | 209 | 210 |
| 211 | 212 | 213 | 214 | 215 | 216 | 217 | 218 | 219 | 220 | 221 | 222 | 223 | 224 | 225 | 226 | 227 | 228 | 229 | 230 | 231 |
| 232 | 233 | 234 | 235 | 236 | 237 | 238 | 239 | 240 | 241 | 242 | 243 | 244 | 245 | 246 | 247 | 248 | 249 | 250 | 251 | 252 |
| 253 | 254 | 255 | 256 | 257 | 258 | 259 | 260 | | | | | | | | | | | | | |
| 274 | 275 | 276 | 277 | 278 | 279 | 280 | 281 | 282 | 283 | 284 | 285 | 286 | 287 | 288 | 289 | 290 | 291 | 292 | 293 | 294 |
| 295 | 296 | 297 | 298 | 299 | 300 | 301 | 302 | 303 | 304 | 305 | 306 | 307 | 308 | 309 | 310 | 311 | 312 | 313 | 314 | 315 |
| 316 | 317 | 318 | 319 | 320 | 321 | 322 | 323 | 324 | 325 | 326 | 327 | 328 | 329 | 330 | 331 | 332 | 333 | 334 | 335 | 336 |
| 337 | 338 | 339 | 340 | 341 | 342 | 343 | 344 | 345 | 346 | 347 | 348 | 349 | 350 | 351 | 352 | 353 | 354 | 355 | 356 | 357 |
| 358 | 359 | 360 | 361 | 362 | 363 | 364 | 365 | 366 | 367 | 368 | 369 | 370 | 371 | 372 | 373 | 374 | 375 | 376 | 377 | 378 |
| 379 | 380 | 381 | 382 | 383 | 384 | 385 | 386 | 387 | 388 | 389 | 390 | 391 | 392 | 393 | 394 | 395 | 396 | 397 | 398 | 399 |
| 400 | 401 | 402 | 403 | 404 | 405 | 406 | 407 | 408 | 409 | 410 | 411 | 412 | 413 | 414 | 415 | 416 | 417 | 418 | 419 | 420 |
| 421 | 422 | 423 | 424 | 425 | 426 | 427 | 428 | 429 | 430 | 431 | 432 | 433 | 434 | 435 | 436 | 437 | 438 | 439 | 440 | 441 |
| 442 | 443 | 444 | 445 | 446 | 447 | 448 | 449 | 450 | 451 | 452 | 453 | 454 | 455 | 456 | 457 | 458 | 459 | 460 | 461 | 462 |
| 463 | 464 | 465 | 466 | 467 | 468 | 469 | 470 | 471 | 472 | 473 | 474 | 475 | 476 | 477 | 478 | 479 | 480 | 481 | 482 | 483 |
| 484 | 485 | 486 | 487 | 488 | 489 | 490 | 491 | 492 | 493 | 494 | 495 | 496 | 497 | 498 | 499 | 500 | 501 | 502 | 503 | 504 |
| 505 | 506 | 507 | 508 | 509 | 510 | 511 | 512 | 513 | 514 | 515 | 516 | 517 | 518 | 519 | 520 | 521 | 522 | 523 | 524 | 525 |
| 526 | 527 | 528 | 529 | 530 | 531 | 532 | 533 | 534 | 535 | 536 | 537 | 538 | 539 | 540 | 541 | 542 | 543 | 544 | 545 | 546 |
| 547 | 548 | 549 | 550 | 551 | 552 | 553 | 554 | 555 | 556 | 557 | 558 | 559 | 560 | 561 | 562 | 563 | 564 | 565 | 566 | 567 |
| 568 | 569 | 570 | 571 | 572 | 573 | 574 | 575 | 576 | 577 | 578 | 579 | 580 | 581 | 582 | 583 | 584 | 585 | 586 | 587 | 588 |
| 589 | 590 | 591 | 592 | 593 | 594 | 595 | 596 | 597 | 598 | 599 | 600 | 601 | 602 | 603 | 604 | 605 | 606 | 607 | 608 | 609 |
| 610 | 611 | 612 | 613 | 614 | 615 | 616 | 617 | 618 | 619 | 620 | 621 | 622 | 623 | 624 | 625 | 626 | 627 | 628 | 629 | 630 |
| 631 | 632 | 633 | 634 | 635 | 636 | 637 | 638 | 639 | 640 | 641 | 642 | 643 | 644 | 645 | 646 | 647 | 648 | 649 | 650 | 651 |

| | | | | |
|---|---|---|---|---|
| white | aliceblue | antiquewhite | antiquewhite1 | antiquewhite2 |
| antiquewhite3 | antiquewhite4 | aquamarine | aquamarine1 | aquamarine2 |
| aquamarine3 | aquamarine4 | azure | azure1 | azure2 |
| azure3 | azure4 | beige | bisque | bisque1 |
| bisque2 | bisque3 | bisque4 | | blanchedalmond |
| blue | blue1 | blue2 | blue3 | blue4 |
| blueviolet | brown | brown1 | brown2 | brown3 |
| brown4 | burlywood | burlywood1 | burlywood2 | burlywood3 |
| burlywood4 | cadetblue | cadetblue1 | cadetblue2 | cadetblue3 |
| cadetblue4 | chartreuse | chartreuse1 | chartreuse2 | chartreuse3 |
| chartreuse4 | chocolate | chocolate1 | chocolate2 | chocolate3 |
| chocolate4 | coral | coral1 | coral2 | coral3 |
| coral4 | cornflowerblue | cornsilk | cornsilk1 | cornsilk2 |
| cornsilk3 | cornsilk4 | cyan | cyan1 | cyan2 |
| cyan3 | cyan4 | darkblue | darkcyan | darkgoldenrod |
| darkgoldenrod1 | darkgoldenrod2 | darkgoldenrod3 | darkgoldenrod4 | darkgray |
| darkgreen | darkgrey | darkkhaki | darkmagenta | darkolivegreen |
| darkolivegreen1 | darkolivegreen2 | darkolivegreen3 | darkolivegreen4 | darkorange |
| darkorange1 | darkorange2 | darkorange3 | darkorange4 | darkorchid |
| darkorchid1 | darkorchid2 | darkorchid3 | darkorchid4 | darkred |
| darksalmon | darkseagreen | darkseagreen1 | darkseagreen2 | darkseagreen3 |
| darkseagreen4 | darkslateblue | darkslategray | darkslategray1 | darkslategray2 |
| darkslategray3 | darkslategray4 | darkslategrey | darkturquoise | darkviolet |
| deeppink | deeppink1 | deeppink2 | deeppink3 | deeppink4 |
| deepskyblue | deepskyblue1 | deepskyblue2 | deepskyblue3 | deepskyblue4 |

**blue = 0**    **blue = 0.2**    **blue = 0.4**

**blue = 0.6**    **blue = 0.8**    **blue = 1**

Quantity of red

Quantity of green

**from Data to Viz**

'From Data to Viz' is a classification of chart types based on input data format. It will help you find the perfect chart in three simple steps :

**1** Identify what type of data you have.

**2** Go to the corresponding decision tree and follow it down to a set of possible charts.

**3** Choose the chart from the set that will suit your data and your needs best.

Dataviz is a world with endless possibilities and this project does not claim to be exhaustive. However it should provide you with a good starting point. For an interactive version and much more, visit:

data-to-viz.com

**WHAT DO YOU WANT TO SHOW ?**

- Distribution
- Correlation
- Ranking
- Part of a whole
- Evolution
- Maps
- Flow

# Data import with the tidyverse : : **CHEATSHEET**

## Read Tabular Data with readr

**read_\***(file, col_names = TRUE, col_types = NULL, col_select = NULL, id = NULL, locale, n_max = Inf, skip = 0, na = c("", "NA"), guess_max = min(1000, n_max), show_col_types = TRUE**)** See **?read_delim**

| A\|B\|C<br>1\|2\|3<br>4\|5\|NA | → | A B C<br>1 2 3<br>4 5 NA | **read_delim(**"file.txt", delim = "\|"**)** Read files with any delimiter. If no delimiter is specified, it will automatically guess.<br>To make file.txt, run: write_file("A\|B\|C\n1\|2\|3\n4\|5\|NA", file = "file.txt") |
|---|---|---|---|
| A,B,C<br>1,2,3<br>4,5,NA | → | A B C<br>1 2 3<br>4 5 NA | **read_csv(**"file.csv"**)** Read a comma delimited file with period decimal marks.<br>write_file("A,B,C\n1,2,3\n4,5,NA", file = "file.csv") |
| A;B;C<br>1,5;2;3<br>4,5;5;NA | → | A B C<br>1.5 2 3<br>4.5 5 NA | **read_csv2(**"file2.csv"**)** Read semicolon delimited files with comma decimal marks.<br>write_file("A;B;C\n1,5;2;3\n4,5;5;NA", file = "file2.csv") |
| A B C<br>1 2 3<br>4 5 NA | → | A B C<br>1 2 3<br>4 5 NA | **read_tsv(**"file.tsv"**)** Read a tab delimited file. Also **read_table()**.<br>**read_fwf(**"file.tsv", fwf_widths(c(2, 2, NA))**)** Read a fixed width file.<br>write_file("A\tB\tC\n1\t2\t3\n4\t5\tNA\n", file = "file.tsv") |

### USEFUL READ ARGUMENTS

**No header**
read_csv("file.csv", col_names = FALSE)

**Provide header**
read_csv("file.csv",
   col_names = c("x", "y", "z"))

**Read multiple files into a single table**
read_csv(c("f1.csv", "f2.csv", "f3.csv"),
   id = "origin_file")

**Skip lines**
read_csv("file.csv", skip = 1)

**Read a subset of lines**
read_csv("file.csv", n_max = 1)

**Read values as missing**
read_csv("file.csv", na = c("1"))

**Specify decimal marks**
read_delim("file2.csv", locale =
   locale(decimal_mark = ";"))

## Save Data with readr

**write_\***(x, file, na = "NA", append, col_names, quote, escape, eol, num_threads, progress**)**

| A B C<br>1 2 3<br>4 5 NA | → | A,B,C<br>1,2,3<br>4,5,NA | |
|---|---|---|---|

**write_delim(**x, file, delim = " "**)** Write files with any delimiter.

**write_csv(**x, file**)** Write a comma delimited file.

**write_csv2(**x, file**)** Write a semicolon delimited file.

**write_tsv(**x, file**)** Write a tab delimited file.

---

One of the first steps of a project is to import outside data into R. Data is often stored in tabular formats, like csv files or spreadsheets.

The front page of this sheet shows how to import and save text files into R using **readr**.

The back page shows how to import spreadsheet data from Excel files using **readxl** or Google Sheets using **googlesheets4**.

## Column Specification with readr

Column specifications define what data type each column of a file will be imported as. By default readr will generate a column spec when a file is read and output a summary.

**spec(**x**)** Extract the full column specification for the given imported data frame.

```
spec(x)
# cols(
#   age = col_integer(),      ← age is an integer
#   edu = col_character(),    ← edu is a character
#   earn = col_double()       ← earn is a double (numeric)
# )
```

### COLUMN TYPES

Each column type has a function and corresponding string abbreviation.

- **col_logical()** - "l"
- **col_integer()** - "i"
- **col_double()** - "d"
- **col_number()** - "n"
- **col_character()** - "c"
- **col_factor(**levels, ordered = FALSE**)** - "f"
- **col_datetime(**format = ""**)** - "T"
- **col_date(**format = ""**)** - "D"
- **col_time(**format = ""**)** - "t"
- **col_skip()** - "-", "_"
- **col_guess()** - "?"

### OTHER TYPES OF DATA

Try one of the following packages to import other types of files:

- **haven** - SPSS, Stata, and SAS files
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)
- **readr::read_lines()** - text data

### USEFUL COLUMN ARGUMENTS

**Hide col spec message**
read_\*(file, show_col_types = FALSE)

**Select columns to import**
Use names, position, or selection helpers.
read_\*(file, col_select = c(age, earn))

**Guess column types**
To guess a column type, read_\*() looks at the first 1000 rows of data. Increase with **guess_max**.
read_\*(file, guess_max = Inf)

### DEFINE COLUMN SPECIFICATION

**Set a default type**
read_csv(
   file,
   col_type = list(.default = col_double())
)

**Use column type or string abbreviation**
read_csv(
   file,
   col_type = list(x = col_double(), y = "l", z = "_")
)

**Use a single string of abbreviations**
# col types: skip, guess, integer, logical, character
read_csv(
   file,
   col_type = "_?ilc"
)

# Import Spreadsheets

## with readxl

### READ EXCEL FILES



**read_excel(**path, sheet = NULL, range = NULL**)**
Read a .xls or .xlsx file based on the file extension.
See front page for more read arguments. Also
**read_xls()** and **read_xlsx()**.
read_excel("excel_file.xlsx")

### READ SHEETS



**read_excel(**path, **sheet = NULL)** Specify which sheet to read by position or name.
read_excel(path, sheet = 1)
read_excel(path, sheet = "s1")

**excel_sheets(**path**)** Get a vector of sheet names.
excel_sheets("excel_file.xlsx")

To **read multiple sheets:**
1. Get a vector of sheet names from the file path.
2. Set the vector names to be the sheet names.
3. Use purrr::map() and purrr::list_rbind() to read multiple files into one data frame.

```
path <- "your_file_path.xlsx"
path |>
  excel_sheets() |>
  set_names() |>
  map(read_excel, path = path) |>
  list_rbind()
```

### OTHER USEFUL EXCEL PACKAGES

For functions to write data to Excel files, see:
- **openxlsx**
- **writexl**

For working with non-tabular Excel data, see:
- **tidyxl**

### READXL COLUMN SPECIFICATION

Column specifications define what data type each column of a file will be imported as.

Use the **col_types** argument of **read_excel()** to set the column specification.

**Guess column types**
To guess a column type, read_excel() looks at the first 1000 rows of data. Increase with the **guess_max** argument.
read_excel(path, guess_max = Inf)

**Set all columns to same type, e.g. character**
read_excel(path, col_types = "text")

**Set each column individually**
```
read_excel(
  path,
  col_types = c("text", "guess", "guess","numeric")
)
```

### COLUMN TYPES

| logical | numeric | text | date | list |
|---------|---------|------|------|------|
| TRUE | 2 | hello | 1947-01-08 | hello |
| FALSE | 3.45 | world | 1956-10-21 | 1 |

- skip
- guess
- logical
- numeric
- text
- date
- list

Use **list** for columns that include multiple data types. See **tidyr** and **purrr** for list-column data.

### CELL SPECIFICATION FOR READXL AND GOOGLESHEETS4



Use the **range** argument of **readxl::read_excel()** or **googlesheets4::read_sheet()** to read a subset of cells from a sheet.
read_excel(path, range = "Sheet1!B1:D2")
read_sheet(ss, range = "B1:D2")

Also use the range argument with cell specification functions **cell_limits()**, **cell_rows()**, **cell_cols()**, and **anchored().**

## with googlesheets4

### READ SHEETS



**read_sheet(**ss, sheet = NULL, range = NULL**)**
Read a sheet from a URL, a Sheet ID, or a dribble from the googledrive package. See front page for more read arguments. Same as **range_read().**

### SHEETS METADATA

**URLs** are in the form:
https://docs.google.com/spreadsheets/d/
**SPREADSHEET_ID**/edit#gid=**SHEET_ID**

**gs4_get(**ss**)** Get spreadsheet meta data.

**gs4_find(**...**)** Get data on all spreadsheet files.

**sheet_properties(**ss**)** Get a tibble of properties for each worksheet. Also **sheet_names()**.

### WRITE SHEETS



**write_sheet(**data, ss = NULL, sheet = NULL**)**
Write a data frame into a new or existing Sheet.

**gs4_create(**name, ..., sheets = NULL**)** Create a new Sheet with a vector of names, a data frame, or a (named) list of data frames.

**sheet_append(**ss, data, sheet = 1**)** Add rows to the end of a worksheet.

### GOOGLESHEETS4 COLUMN SPECIFICATION

Column specifications define what data type each column of a file will be imported as.

Use the **col_types** argument of **read_sheet()/ range_read()** to set the column specification.

**Guess column types**
To guess a column type read_sheet()/ range_read() looks at the first 1000 rows of data. Increase with **guess_max**.
read_sheet(path, guess_max = Inf)

**Set all columns to same type, e.g. character**
read_sheet(path, col_types = "c")

**Set each column individually**
# col types: skip, guess, integer, logical, character
read_sheets(ss, col_types = "_?ilc")

### COLUMN TYPES

| I | n | c | D | L |
|------|------|-------|------------|-------|
| TRUE | 2 | hello | 1947-01-08 | hello |
| FALSE | 3.45 | world | 1956-10-21 | 1 |

- skip - "_" or "-"
- guess - "?"
- logical - "l"
- integer - "i"
- double - "d"
- numeric - "n"
- date - "D"
- datetime - "T"
- character - "c"
- list-column - "L"
- cell - "C" Returns list of raw cell data.

Use list for columns that include multiple data types. See **tidyr** and **purrr** for list-column data.

### FILE LEVEL OPERATIONS

**googlesheets4** also offers ways to modify other aspects of Sheets (e.g. freeze rows, set column width, manage (work)sheets). Go to **googlesheets4.tidyverse.org** to read more.
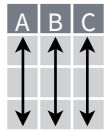
For whole-file operations (e.g. renaming, sharing, placing within a folder), see the tidyverse package **googledrive** at **googledrive.tidyverse.org**.
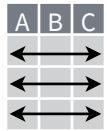
# Data tidying with tidyr : : **CHEATSHEET**

**Tidy data** is a way to organize tabular data in a consistent data structure across packages. A table is tidy if:
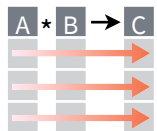
Each **variable** is in its own **column** & Each **observation**, or **case**, is in its own row

Access **variables** as **vectors**

Preserve **cases** in vectorized operations

## Tibbles

### AN ENHANCED DATA FRAME

Tibbles are a table format provided by the **tibble** package. They inherit the data frame class, but have improved behaviors:

- **Subset** a new tibble with ], a vector with [[ and $.
- **No partial matching** when subsetting columns.
- **Display** concise views of the data on one screen.

**options**(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf**)** Control default display settings.

**View()** or **glimpse()** View the entire data set.

### CONSTRUCT A TIBBLE

**tibble(**…**)** Construct by columns.
tibble(x = 1:3, y = c("a", "b", "c"))

**tribble(**…**)** Construct by rows.
tribble(~x,  ~y,
        1, "a",
        2, "b",
        3, "c")

**Both make this tibble**
```
A tibble: 3 × 2
      x     y
  <int> <chr>
1     1     a
2     2     b
3     3     c
```

**as_tibble(**x, …**)** Convert a data frame to a tibble.

**enframe(**x, name = "name", value = "value"**)** Convert a named vector to a tibble. Also **deframe()**.

**is_tibble(**x**)** Test whether x is a tibble.

## Reshape Data - Pivot data to reorganize values into a new layout.

**pivot_longer(**data, cols, names_to = "name", values_to = "value", values_drop_na = FALSE**)**

"Lengthen" data by collapsing several columns into two. Column names move to a new names_to column and values to a new values_to column.

pivot_longer(table4a, cols = 2:3, names_to ="year", values_to = "cases")

**pivot_wider(**data, names_from = "name", values_from = "value"**)**

The inverse of pivot_longer(). "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

pivot_wider(table2, names_from = type, values_from = count)

## Split Cells - Use these functions to split or combine cells into individual, isolated values.

**unite(**data, col, …, sep = "_", remove = TRUE, na.rm = FALSE**)** Collapse cells across several columns into a single column.

unite(table5, century, year, col = "year", sep = "")

**separate_wider_delim(**data, cols, delim, …, names = NULL, names_sep = NULL, names_repair = "check unique", too_few, too_many, cols_remove = TRUE**)** Separate each cell in a column into several columns. Also **separate_wider_regex()** and **separate_wider_position().**

separate(table3, rate, sep = "/",
    into = c("cases", "pop"))

**separate_longer_delim(**data, cols, delim, .., width, keep_eampty**)** Separate each cell in a column into several rows.

separate_longer_delim(table3, rate, sep = "/")

## Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

**expand(**data, …**)** Create a new tibble with all possible combinations of the values of the variables listed in … Drop other variables.
expand(mtcars, cyl, gear, carb)

**complete(**data, …, fill = list()**)** Add missing possible combinations of values of variables listed in … Fill remaining variables with NA.
complete(mtcars, cyl, gear, carb)

## Handle Missing Values

Drop or replace explicit missing values (NA).

**drop_na(**data, …**)** Drop rows containing NA's in … columns.
drop_na(x, x2)

**fill(**data, …, .direction = "down"**)** Fill in NA's in … columns using the next or previous value.
fill(x, x2)

**replace_na(**data, replace**)** Specify a value to replace NA in selected columns.
replace_na(x, list(x2 = 2))

# Nested Data

A **nested data frame** stores individual tables as a list-column of data frames within a larger organizing data frame. List-columns can also be lists of vectors or lists of varying data types. Use a nested data frame to:
- Preserve relationships between observations and subsets of data. Preserve the type of the variables being nested (factors and datetimes aren't coerced to character).
- Manipulate many sub-tables at once with **purrr** functions like map(), map2(), or pmap() or with **dplyr** rowwise() grouping.

## CREATE NESTED DATA

**nest(data, …)** Moves groups of cells into a list-column of a data frame. Use alone or with dplyr::**group_by()**:

1. Group the data frame with **group_by()** and use **nest()** to move the groups into a list-column.
   ```
   n_storms <- storms |>
     group_by(name) |>
     nest()
   ```

2. Use **nest(new_col = c(x, y))** to specify the columns to group using dplyr::**select()** syntax.
   ```
   n_storms <- storms |>
     nest(data = c(year:long))
   ```



Index list-columns with [[]]. n_storms$data[[1]]

## CREATE TIBBLES WITH LIST-COLUMNS

tibble::**tribble(…)** Makes list-columns when needed.
```
tribble( ~max, ~seq,
          3,   1:3,
          4,   1:4,
          5,   1:5)
```

| max | seq |
|-----|-----|
| 3 | <int [3]> |
| 4 | <int [4]> |
| 5 | <int [5]> |

tibble::**tibble(…)** Saves list input as list-columns.
```
tibble(max = c(3, 4, 5), seq = list(1:3, 1:4, 1:5))
```

tibble::**enframe(**x, name="name", value="value"**)**
Converts multi-level list to a tibble with list-cols.
```
enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')
```

## OUTPUT LIST-COLUMNS FROM OTHER FUNCTIONS

dplyr::**mutate()**, **transmute()**, and **summarise()** will output list-columns if they return a list.
```
mtcars |>
  group_by(cyl) |>
  summarise(q = list(quantile(mpg)))
```

## RESHAPE NESTED DATA

**unnest(**data, cols, …, keep_empty = FALSE**)** Flatten nested columns back to regular columns. The inverse of nest().
n_storms |> unnest(data)

**unnest_longer(**data, col, values_to = NULL, indices_to = NULL**)**
Turn each element of a list-column into a row.
```
starwars |>
  select(name, films) |>
  unnest_longer(films)
```



**unnest_wider(**data, col**)** Turn each element of a list-column into a regular column.
```
starwars |>
  select(name, films) |>
  unnest_wider(films, names_sep = "_")
```



**hoist(**.data, .col, …, .remove = TRUE**)** Selectively pull list components out into their own top-level columns. Uses purrr::pluck() syntax for selecting from lists.
```
starwars |>
  select(name, films) |>
  hoist(films, first_film = 1, second_film = 2)
```



## TRANSFORM NESTED DATA

A vectorized function takes a vector, transforms each element in parallel, and returns a vector of the same length. By themselves vectorized functions cannot work with lists, such as list-columns.

dplyr::**rowwise(**.data, …**)** Group data so that each row is one group, and within the groups, elements of list-columns appear directly (accessed with [[ ), not as lists of length one. **When you use rowwise(), dplyr functions will seem to apply functions to list-columns in a vectorized fashion.**



Apply a function to a list-column and **create a new list-column.**
```
n_storms |>
  rowwise() |>
  mutate(n = list(dim(data)))
```
> dim() returns two values per row
> wrap with list to tell mutate to create a list-column

Apply a function to a list-column and **create a regular column.**
```
n_storms |>
  rowwise() |>
  mutate(n = nrow(data))
```
> nrow() returns one integer per row

Collapse **multiple list-columns** into a single list-column.
```
starwars |>
  rowwise() |>
  mutate(transport = list(append(vehicles, starships)))
```
> append() returns a list for each row, so col type must be list

Apply a function to **multiple list-columns.**
```
starwars |>
  rowwise() |>
  mutate(n_transports = length(c(vehicles, starships)))
```
> length() returns one integer per row

See **purrr** package for more list functions.

# R GRAPHICAL PARAMETERS CHEATSHEET

## Box feature

**bty** : kind of box

*o=complete / ?=top & right / n=no box / c=top & left & bottom / l=bottom & left*

## Title

**main** : name of the title
**cex.main** : size *cex.main=2*
**col.main** : color *col.main="red"*
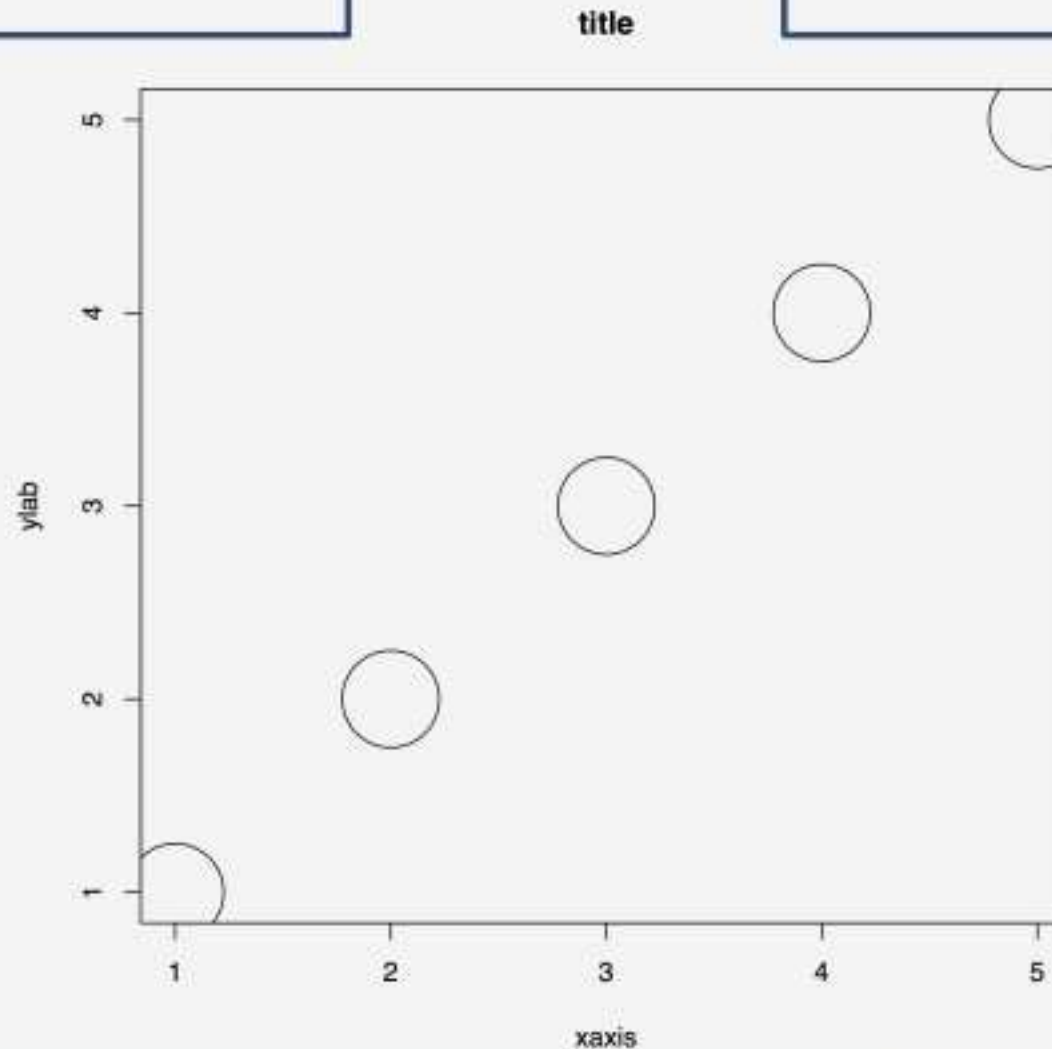**font.main** : font *font.main=3*

title

## Margins

See **Graph #135**
*mar, oma, omd, omi*

## General

**tck** : add a grid on a % of the area *tck=1*
**bg** : color of the background *par(bg="red")*
**font** : font of the text *(normal, bold, italic..)*
**lheight**: size between lines of titles
**srt**: text rotation in degree

ylab

## Symbol styles

See **Graph #6**
*pch, lwd, pty, col, cex, type...*

xaxis

## X-axis name

**xlab** : name of the axis
**cex.lab** : size *cex.lab=2*
**col.lab** : color *col.lab="red"*
**sub** : to add a subtitle

## X-axis features

**lab** : number of graduation *lab=c(12,2,0)*
**xaxp** : to add c graduation from a to b: *xaxp=c(a,b,c)*
**log** : for logarithmic scale: *log="x"*
**xaxt** : to remove x axis: *xaxt="n"*
**fg** : color of axis, ticks and grid: *fg="red"*
**cex.axis** : size of tick labels
**col.axis** : color of tick labels
**xlim** : limits of the axis *xlim=c(0,10)*
**las**: orientation of tick labels
*0=parralel to the axis, 1=horizontal...*

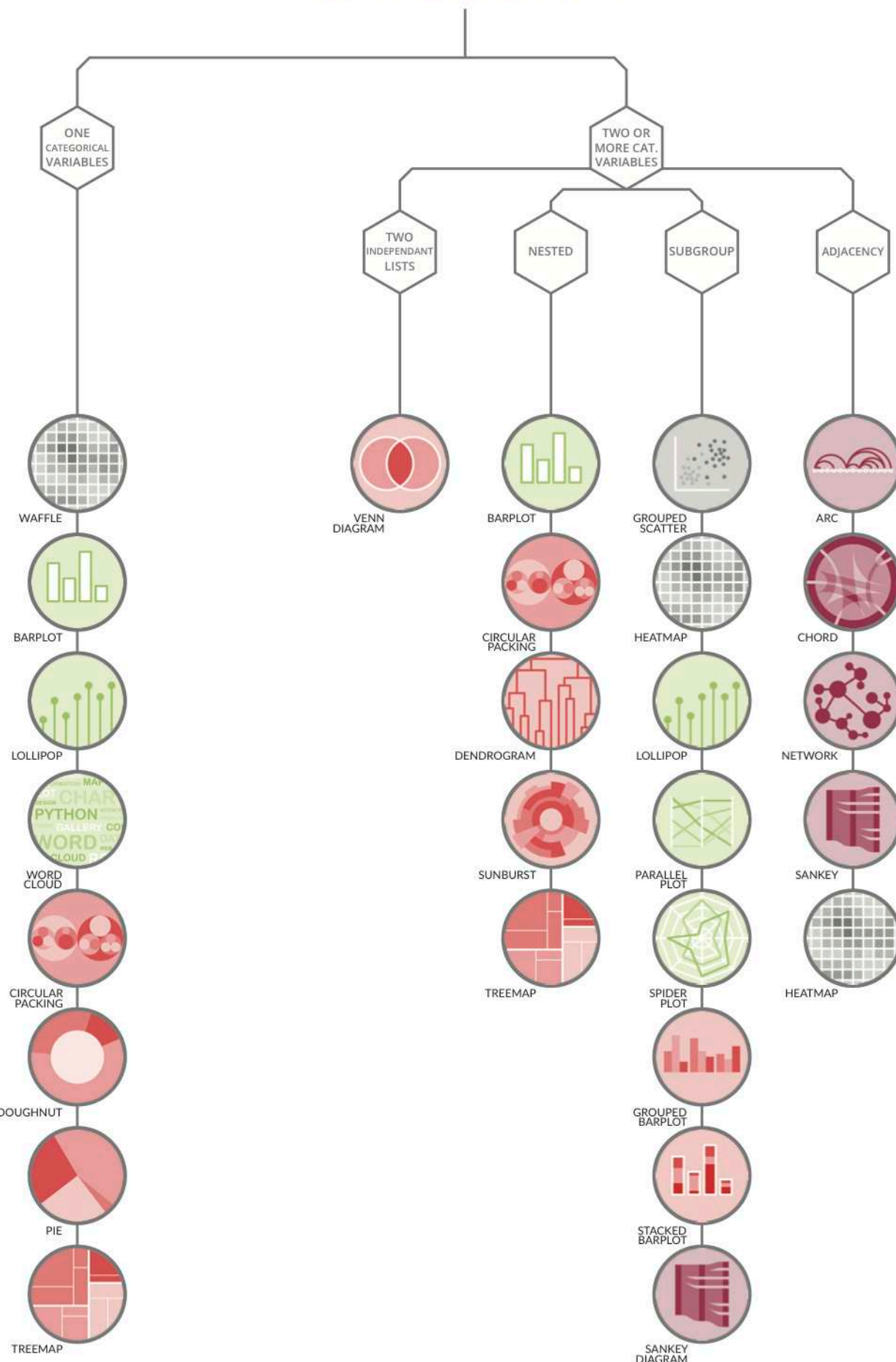| | | | | |
|---|---|---|---|---|
| white | aliceblue | antiquewhite | antiquewhite1 | antiquewhite2 |
| antiquewhite3 | antiquewhite4 | aquamarine | aquamarine1 | aquamarine2 |
| aquamarine3 | aquamarine4 | azure | azure1 | azure2 |
| azure3 | azure4 | beige | bisque | bisque1 |
| bisque2 | bisque3 | bisque4 | | blanchedalmond |
| blue | blue1 | blue2 | blue3 | blue4 |
| blueviolet | brown | brown1 | brown2 | brown3 |
| brown4 | burlywood | burlywood1 | burlywood2 | burlywood3 |
| burlywood4 | cadetblue | cadetblue1 | cadetblue2 | cadetblue3 |
| cadetblue4 | chartreuse | chartreuse1 | chartreuse2 | chartreuse3 |
| chartreuse4 | chocolate | chocolate1 | chocolate2 | chocolate3 |
| chocolate4 | coral | coral1 | coral2 | coral3 |
| coral4 | cornflowerblue | cornsilk | cornsilk1 | cornsilk2 |
| cornsilk3 | cornsilk4 | cyan | cyan1 | cyan2 |
| cyan3 | cyan4 | darkblue | darkcyan | darkgoldenrod |
| darkgoldenrod1 | darkgoldenrod2 | darkgoldenrod3 | darkgoldenrod4 | darkgray |
| darkgreen | darkgrey | darkkhaki | darkmagenta | darkolivegreen |
| darkolivegreen1 | darkolivegreen2 | darkolivegreen3 | darkolivegreen4 | darkorange |
| darkorange1 | darkorange2 | darkorange3 | darkorange4 | darkorchid |
| darkorchid1 | darkorchid2 | darkorchid3 | darkorchid4 | darkred |
| darksalmon | darkseagreen | darkseagreen1 | darkseagreen2 | darkseagreen3 |
| darkseagreen4 | darkslateblue | darkslategray | darkslategray1 | darkslategray2 |
| darkslategray3 | darkslategray4 | darkslategrey | darkturquoise | darkviolet |
| deeppink | deeppink1 | deeppink2 | deeppink3 | deeppink4 |
| deepskyblue | deepskyblue1 | deepskyblue2 | deepskyblue3 | deepskyblue4 |

from Data to Viz

**'From Data to Viz'** is a classification of chart types based on input data format. It will help you find the perfect chart in three simple steps :

(1) Identify what type of data you have.

(2) Go to the corresponding decision tree and follow it down to a set of possible charts.

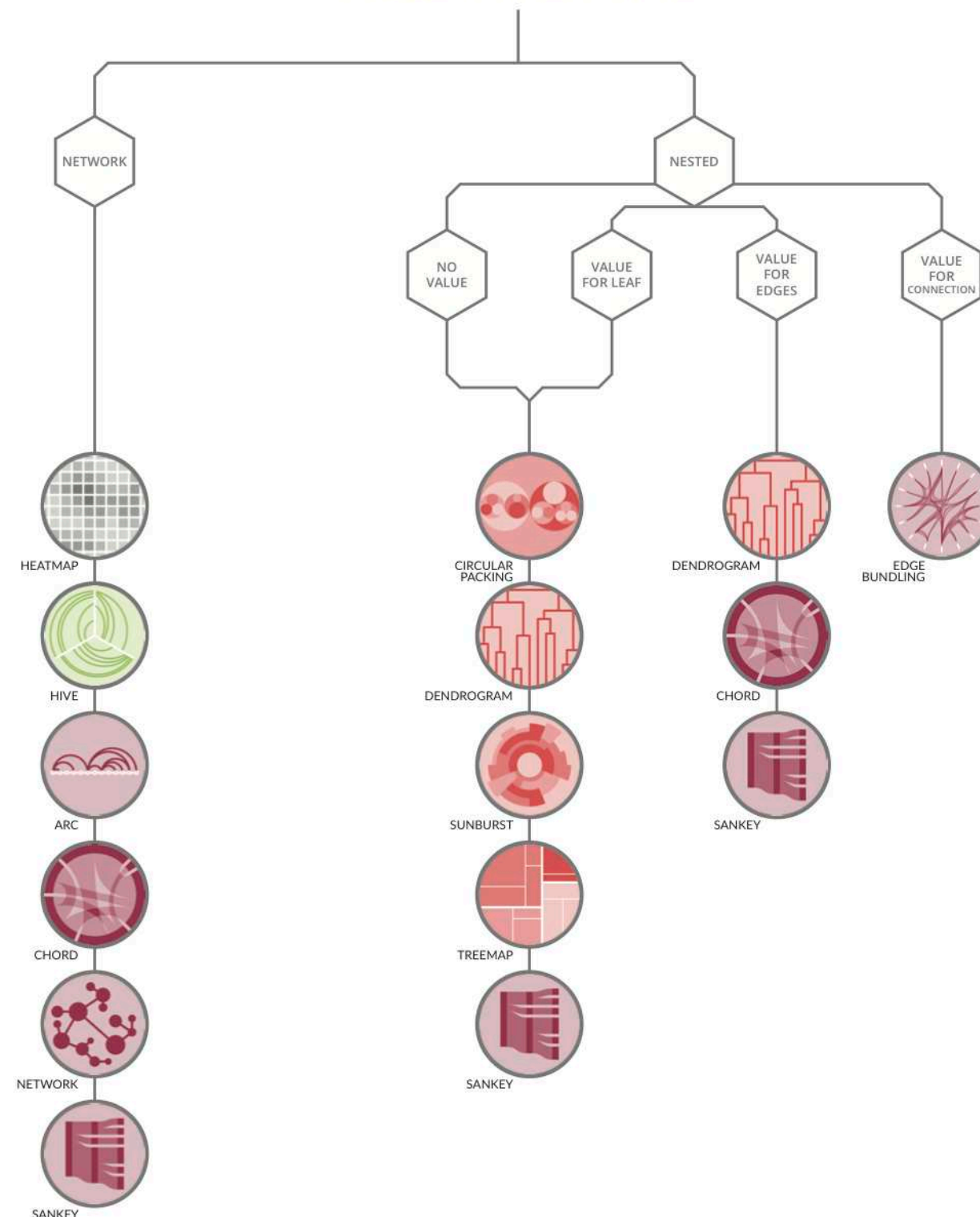(3) Choose the chart from the set that will suit your data and your needs best.

Dataviz is a world with endless possibilities and this project does not claim to be exhaustive. However it should provide you with a good starting point. For an interactive version and much more, visit:
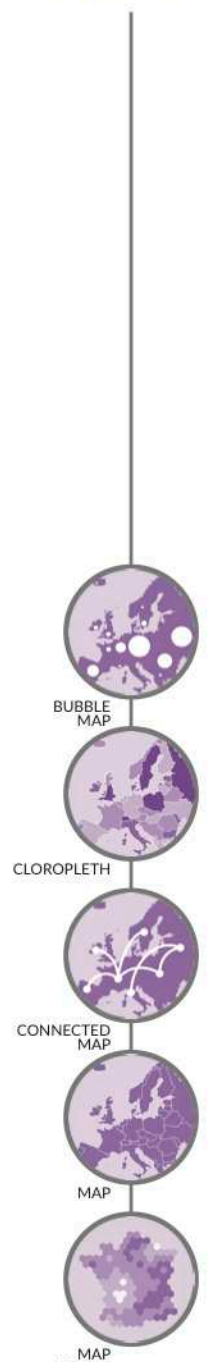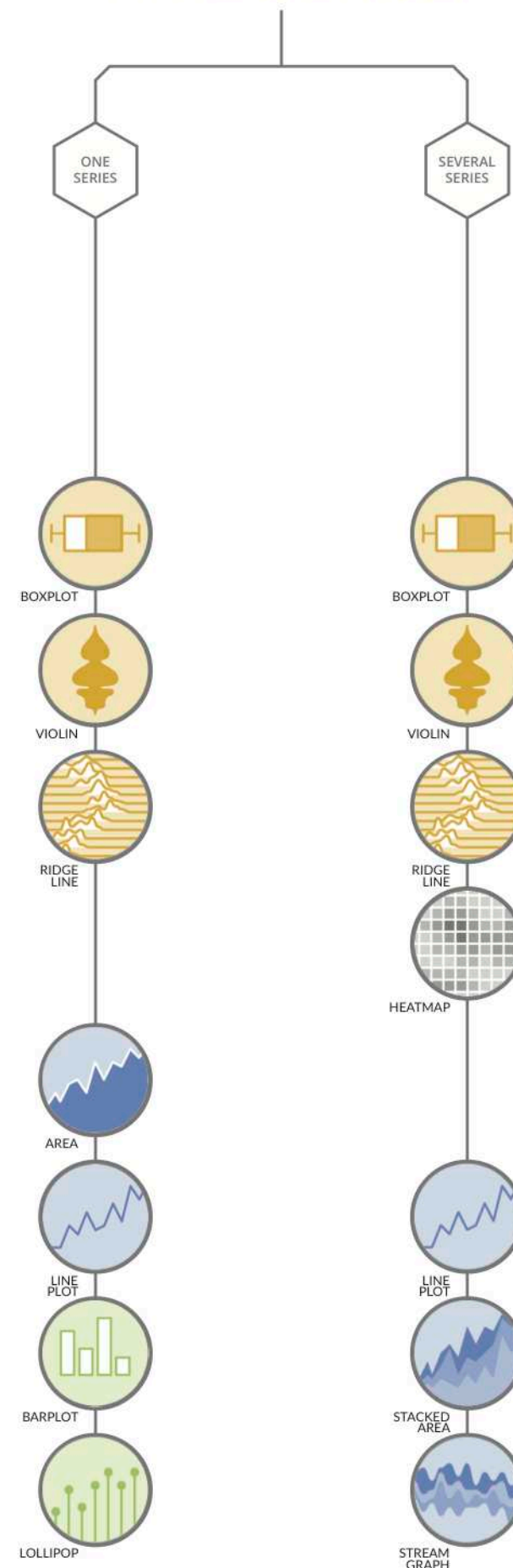
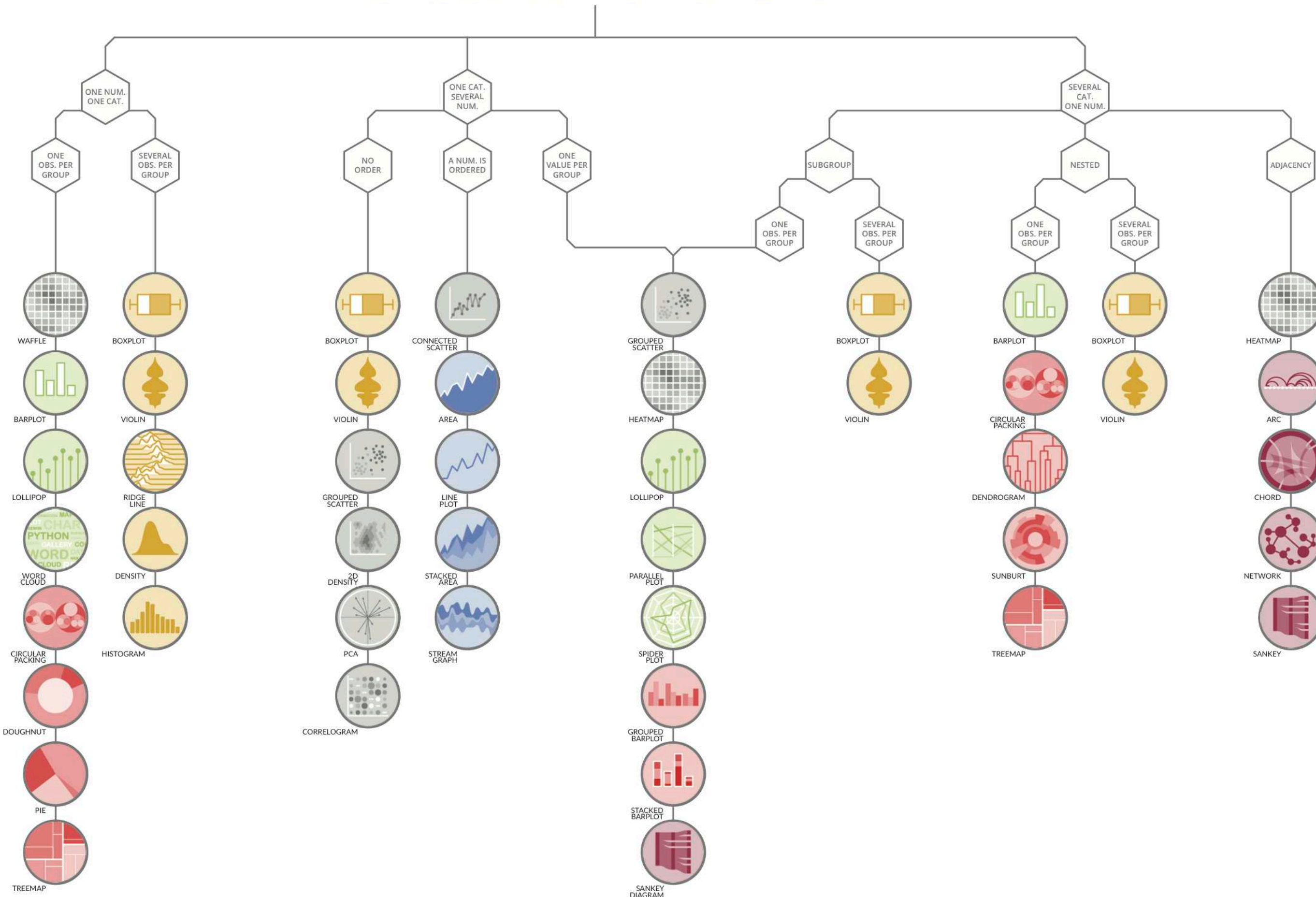data-to-viz.com

## CATEGORIC
## RELATIONAL
## MAP
## TIME SERIES
## CATEGORIC AND NUMERIC
## NUMERIC

WHAT DO YOU WANT TO SHOW ?

- Distribution
- Correlation
- Ranking
- Part of a whole
- Evolution
- Maps
- Flow

2018 © Yan Holtz & Conor Healy for www.data-to-viz.com