

# Midterm I

April 20 2023

Name:

**Total Points: 100**

*Gapminder data*

Health and income outcomes for 142 countries from 1952 to 2007 in increments of 5 years. The variables in the dataset are `country`, `continent`, `year`, `lifeExp`, `pop`, and `gdpPercap`. The descriptions for the variables are:

- `country` : name of the country, factor with 142 levels
- `continent`: name of the continent, factor with 5 levels
- `year` : ranges from 1952 to 2007 in increments of 5 years (12 distinct years)
- `lifeExp`: life expectancy at birth, in years
- `pop` : population
- `gdpPercap` : GDP per capita (US\$, inflation-adjusted)

```
glimpse(gapminder)
Rows: 1,704
Columns: 6
$ country    <fct> "Afghanistan", "Afghanistan", "Afghanistan", "Afghanistan", ~
$ continent <fct> Asia, Asia, Asia, Asia, Asia, Asia, Asia, Asia, Asia, ~
$ year       <int> 1952, 1957, 1962, 1967, 1972, 1977, 1982, 1987, 1992, 1997, ~
$ lifeExp    <dbl> 28.801, 30.332, 31.997, 34.020, 36.088, 38.438, 39.854, 40.8~
$ pop        <int> 8425333, 9240934, 10267083, 11537966, 13079460, 14880372, 12~
$ gdpPercap   <dbl> 779.4453, 820.8530, 853.1007, 836.1971, 739.9811, 786.1134, ~
```

```
gapminder %>% pull(continent) %>% unique()
[1] Asia      Europe    Africa    Americas Oceania
Levels: Africa Americas Asia Europe Oceania
```

## Part 1: Data Wrangling (*10 points each*)

**What do the following code chunks do?** Provide a thorough and intuitive (3-5 sentences) description of the output from each of the following R chunks. The chunks produce a new data set. Please give the dimensions in addition to your description. Write your descriptions in regular English, without using variable names.

a.

```
gapminder %>%
  filter(year %in% c(1952, 2007)) %>%
  group_by(continent, year) %>%
  summarize(median_gdpPercap = median(gdpPercap),
            median_lifeExp = median(lifeExp),
            median_pop = median(pop))
```

b.

```
gapminder %>%
  mutate(gdp_total = gdpPercap * pop) %>%
  group_by(continent, year) %>%
  summarize(gdp_continent = sum(gdp_total)) %>%
  pivot_wider(names_from = year,
              values_from = gdp_continent,
              names_prefix = "year_")
```

c.

```
set.seed(143)
selected_countries <- gapminder %>%
  distinct(country, continent) %>%
  group_by(continent) %>%
  slice_sample(n = 1) %>%
  pull(country)

gapminder %>%
  filter(country %in% selected_countries) %>%
  filter(year %in% c(1952, 2007)) %>%
  group_by(country, year) %>%
  summarize(median_gdpPercap = median(gdpPercap),
            median_lifeExp = median(lifeExp),
            total_pop = sum(pop)) %>%
  pivot_longer(cols = -c(country, year),
               names_to = "stat",
               values_to = "value")
```

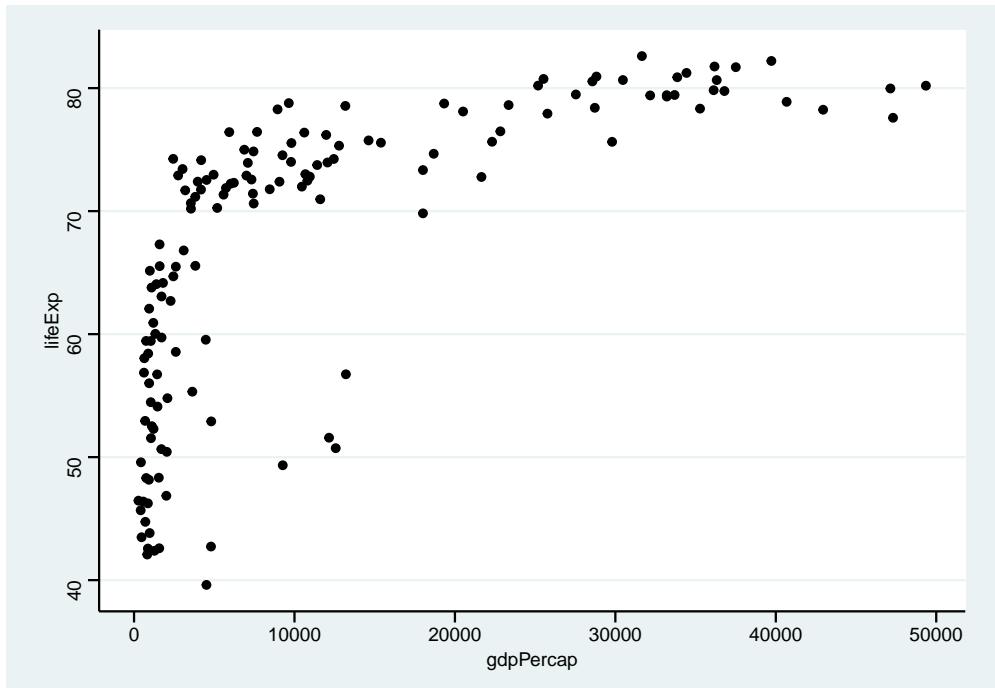
d.

```
set.seed(143)
gapminder %>%
  filter(continent == "Asia") %>%
  distinct(country) %>%
  slice_sample(n = 5) %>%
  inner_join(gapminder, by = "country") %>%
  filter(year %in% c(1952, 2007)) %>%
  group_by(country, year) %>%
  summarize(median_gdpPercap = median(gdpPercap),
            median_lifeExp = median(lifeExp),
            total_pop = sum(pop)) %>%
  pivot_longer(cols = -c(country, year),
               names_to = "stat",
               values_to = "value") %>%
  pivot_wider(names_from = year,
              values_from = value,
              names_prefix = "year_")
```

## Part 2: Graphics (*15 points each*)

- a. The scatter plot below visualizes the relationship between GDP per capita and life expectancy of countries in the year 2007. What are 5 ways you could improve the aesthetics and readability of this plot by following best data visualization practices? Also write 5 code modifications on the space provided below:

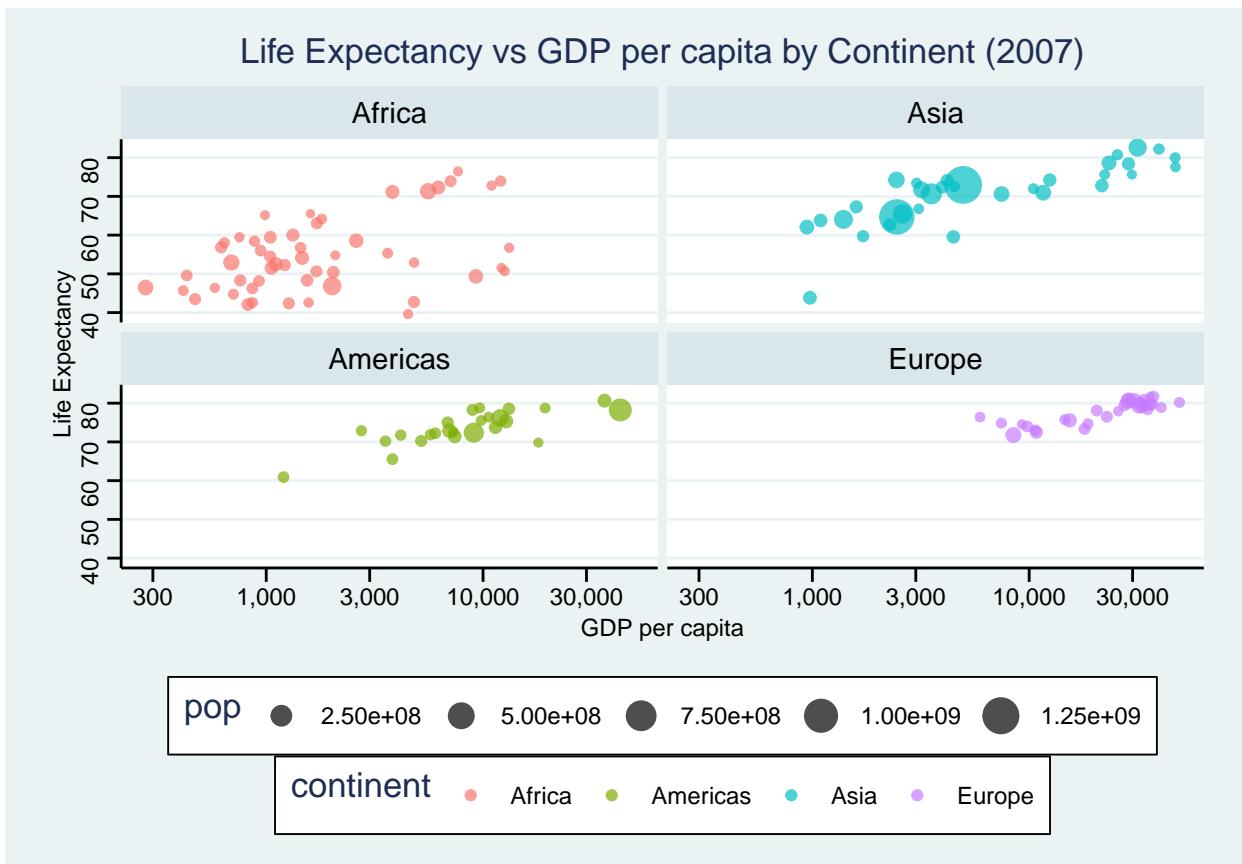
```
gapminder %>%
  filter(year == 2007) %>%
  ggplot(aes(x = gdpPercap, y = lifeExp)) +
  geom_point()
```



*Answer:*

```
gapminder %>%
  filter(year == 2007) %>%
  ggplot(aes(x = gdpPercap, y = lifeExp)) +
```

b. The partial code used to generate the plot below is given with placeholders for code snippet. Please provide the appropriate code snippet.



```
gapminder %>%
  ##### FILL IN i. ##### %>%
  ##### FILL IN ii. ##### %>%
  group_by(continent) %>%
  mutate(avg_lifeExp = mean(lifeExp)) %>%
  ungroup() %>%
  ##### FILL IN iii. ##### %>%
  ##### FILL IN iv. ##### +
  geom_point(alpha = 0.7) +
  labs(title = "Life Expectancy vs GDP per capita by Continent (2007)",
       x = "GDP per capita",
       y = "Life Expectancy") +
  scale_fill_brewer(palette = "Set1") +
  scale_x_log10(labels = scales::comma) +
  ##### FILL IN v. ##### +
  theme(legend.position = "bottom")
```

- i. Create a new column called “year\_date” that converts the “year” column into a date object
- ii. Filter the data to only include rows where the year is 2007 and exclude continent “Oceania”
- iii. Reorder the levels of the “continent” factor based on the “avg\_lifeExp” column and store the result in a new column called “continent\_reordered.”
- iv. Create a ggplot2 scatter plot with “gdpPercap” on the x-axis, “lifeExp” on the y-axis, point sizes representing the “pop” column, and points color-coded by the “continent” column.
- v. Create a faceted plot based on the “continent\_reordered” column, with 2 columns of panes.

### Part 3: Data Objects (*5 points each*)

```
x <- 4:1
y <- c(TRUE, factor(c(NA, "b")), 1)
z <- list(z1 = x, z2 = y, z3 = c("Carleton", "college"), z4 = matrix(1:9, nrow = 3))
```

Consider the above objects to tell what each of the following code chunks evaluate to? Briefly explain your answer.

(a)

```
y[3]
```

(b)

```
z[["z3"]][1]
```

(c)

```
z[x][[2]][[2]]
```

(d)

```
x - y
```

(e)

```
unlist(z)
```

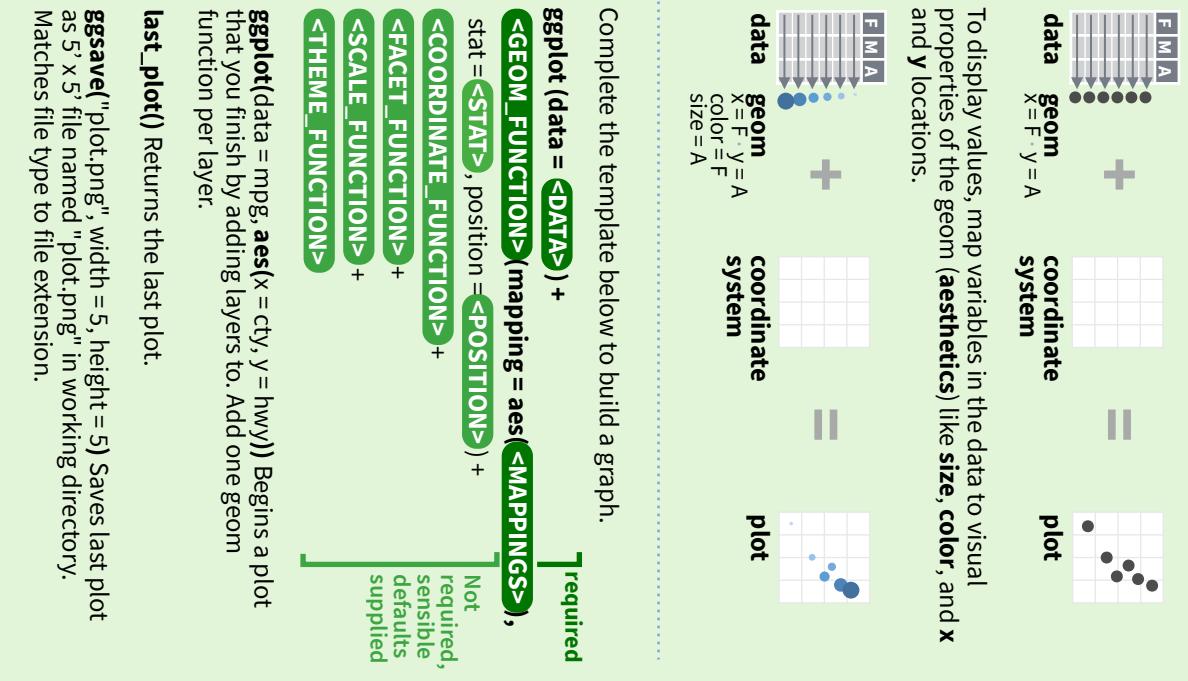
(f)

```
typeof(unlist(z))
```

# Data visualization with ggplot2 :: CHEAT SHEET

## 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.



## 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))
```

### TWO VARIABLES both continuous

**e + geom\_label()** **continuous bivariate distribution**  
`e + geom_label(aes(label = cty), nudge_x = 1, nudge_y = 1) -> x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust`

**B**

**C**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

**T**

**U**

**V**

**W**

**X**

**Y**

**Z**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

**T**

**U**

**V**

**W**

**X**

**Y**

**Z**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

**T**

**U**

**V**

**W**

**X**

**Y**

**Z**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

**T**

**U**

**V**

**W**

**X**

**Y**

**Z**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

**T**

**U**

**V**

**W**

**X**

**Y**

**Z**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

**T**

**U**

**V**

**W**

**X**

**Y**

**Z**

**A**

**B**

**C**

**D**

**E**

**F**

**G**

**H**

**I**

**J**

**K**

**L**

**M**

**N**

**O**

**P**

**Q**

**R**

**S**

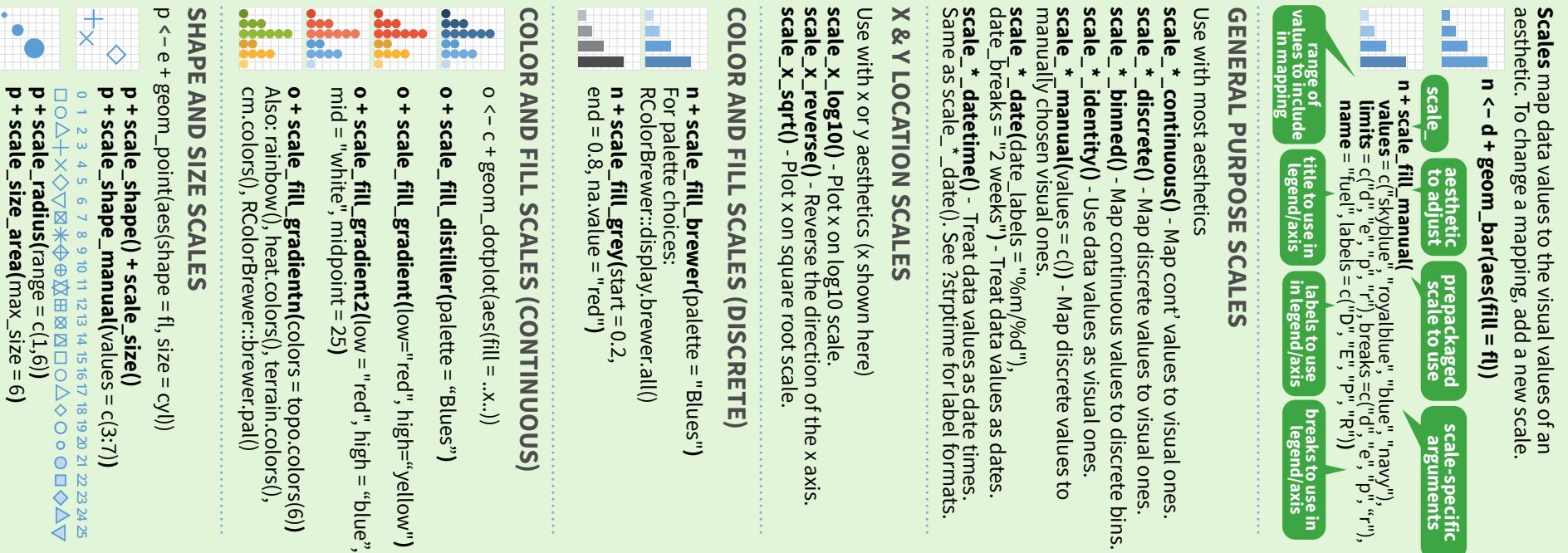


Stats

An alternative way to build a layer.



Scales



## Coordinate Systems



Faceting

1  
-  
-  
-  
-



# Data transformation with dplyr :: CHEAT SHEET

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



**pipes**

Each **variable** is in its own **column**. Each **observation**, or **case**, is in its own **row**.

**x %>% f(y)** becomes **f(x, y)**

## Summarise 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**

**summarise(data, ...)**  
Compute table of summaries.

summarise(mtcars, 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()**.  
count(mtcars, cyl)

## Group Cases

Use **group\_by(data, ...)** 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) %>% summarise(avg = mean(mpg))**

**mtcars %>% group\_by(cyl) %>% summarise(avg = mean(mpg))**

## ARRANGE CASES

Use **rowwise(data, ...)** to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See [tidy cheatsheet](#) for list-column workflow.

**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.  
arrange(mtcars, mpg)  
arrange(mtcars, desc(mpg))

## ADD CASES

**ungroup(x, ...)** Returns ungrouped copy of table.  
ungroup(g\_mtcars)

**add\_row(data, ..., before = NULL, after = NULL)** Add one or more rows to a table.  
add\_row(cars, speed = 1, dist = 1)

**transmute(data, ...)** Compute new column(s), drop others.  
transmute(mtcars, gpm = 1 / mpg)

**rename(data, ...)** Rename columns. Use **rename\_with()** to rename with a function.  
rename(cars, distance = dist)

## 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.  
filter(mtcars, mpg > 20)

**distinct(data, ..., .keep\_all = FALSE)** Remove rows with duplicate values.  
distinct(mtcars, gear)

**slice(data, ..., .preserve = FALSE)** Select rows by position.  
slice(mtcars, 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.  
slice\_sample(mtcars, 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.  
slice\_min(mtcars, mpg, prop = 0.25)

**slice\_head(data, ..., n, prop)** and **slice\_tail()** Select the first or last rows.  
slice\_head(mtcars, n = 5)

**slice\_head(data, ..., n, prop)** and **slice\_tail()** Select the first or last rows.  
slice\_head(mtcars, n = 5)

**slice\_head(data, ..., n, prop)** and **slice\_tail()** Select the first or last rows.  
slice\_head(mtcars, n = 5)

### MANIPULATE MULTIPLE VARIABLES AT ONCE

**across(cols, funs, ...)** Summarise or mutate multiple columns in the same way.  
summarise(mtcars, across(everything(), mean))

**c\_across(.cols)** Compute across columns in row-wise data.  
transmute(rowwise(UKgas), 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**

**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.  
arrange(mtcars, mpg)  
arrange(mtcars, desc(mpg))

**mutate(data, ..., .keep = "all", before = NULL, after = NULL)** Compute new column(s). Also **add\_column()**, **add\_count()**, and **add\_tally()**.  
mutate(mtcars, gpm = 1 / mpg)

## 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.  
pull(mtcars, wt)

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

**relocate(data, ..., .before = NULL, .after = NULL)** Move columns to new position.  
relocate(mtcars, mpg, cyl, .after = last\_col())

**select\_helpers** with **select()** and **across()**  
e.g. **select(mtcars, mpg:cyl)**

**contains(match)** **num\_range(prefix, range)** ; e.g. mpg:cyl  
**ends\_with(match)** **all\_of(x)/any\_of(x, ...)** ; e.g. -gear  
**starts\_with(match)** **matches(match)** **everything()**



# Vectorized Functions

## TO USE WITH MUTATE()

**mutate()** and **transmute()** apply 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



dplyr::lag() - offset elements by 1  
dplyr::lead() - offset elements by -1

## CUMULATIVE AGGREGATE

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

## RANKING

dplyr::cume\_dist() - proportion of all values <= rank w ties = min, no gaps  
dplyr::dense\_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 ~ "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 SUMMARISE()

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

## summary function



dplyr::n() - number of values/rows  
dplyr::n\_distinct() - # of uniques  
sum(!is.na()) - # of non-NA's

## COUNT

dplyr::n() - number of values/rows  
dplyr::n\_distinct() - # of uniques  
sum(!is.na()) - # of non-NA's

## POSITION

mean() - mean, also mean(!is.na())  
median() - median

## LOGICAL

mean() - proportion of TRUE's  
sum() - # of TRUE's

## 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 <- rownames_to_column(mtcars,  
#> #> var = "C")
```

```
tibble::column_to_rownames()  
#> #> Move col into row names.  
#> #> column_to_rownames(a, 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"), ...)** Join matching na\_matches = "na" Join matching values from y to x.

**right\_join(x, y, by = NULL, copy = FALSE, suffix = c(".X", ".Y"), ...)** Keep matching na\_matches = "na" Join matching values from x to y.

**inner\_join(x, y, by = NULL, copy = 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.

## 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).



## COMBINE VARIABLES

**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).

## SET OPERATIONS

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

**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

**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.

**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.

**left\_join(x, y, by = c("C" = "D"))** left\_join(x, y, by = c("C" = "D"))

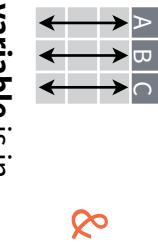
**use\_suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.

**use\_setequal** to test whether two data sets contain the exact same rows (in any order).

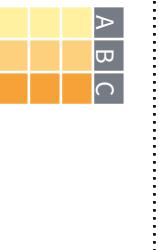
# Data tidying With `tidyverse`

**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 **case**, is in its own row



Each **observation**, or its own **column**



## 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.

**tribble(...)** Construct by rows.

**Both make this tibble**



**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.

| country | 1999 | 2000 |
|---------|------|------|
| A       | 0.7K | 2K   |
| B       | 37K  | 80K  |
| C       | 212K | 213K |
| A       | 1999 | 2K   |
| B       | 2000 | 80K  |
| C       | 2000 | 213K |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

| country | year | cases |
| --- | --- | --- |
| A | 1999 | 0.7K |
| B | 1999 | 37K |
<tbl\_info cols

