

“Data verbs” cheatsheet

Each column is a single **named variable**

These rules are also known as: **keeping data “tidy”**

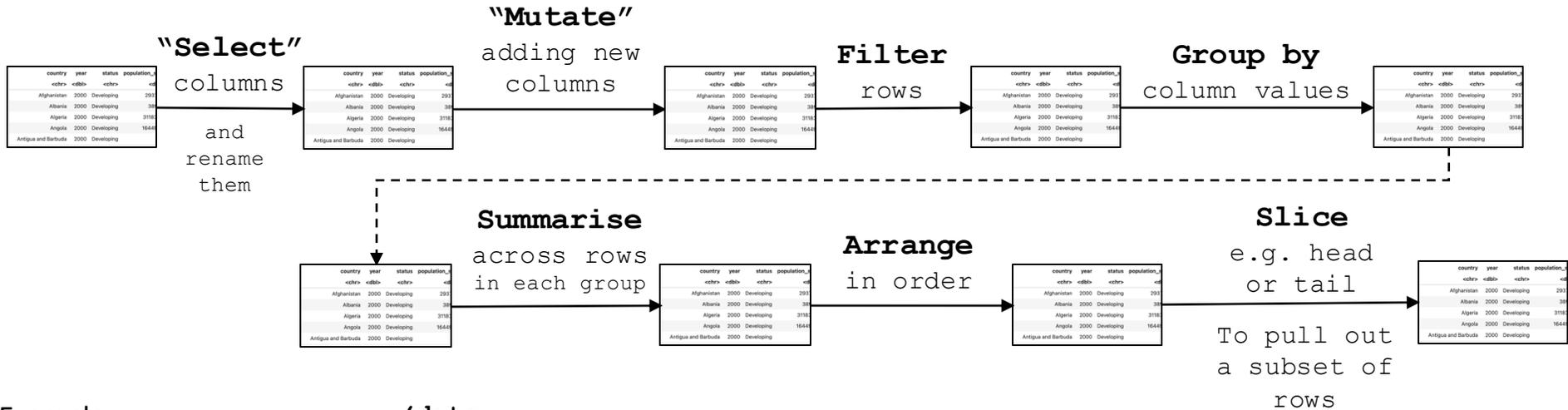
1. Keep your data in a **data frame**

data =

country	year	status	population_size	life_expectancy	adult_mortality	under_five_deaths
<chr>	<dbl>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
Afghanistan	2000	Developing	293756	54.8	321	122
Albania	2000	Developing	38927	72.6	11	1
Algeria	2000	Developing	3118366	71.3	145	25
Angola	2000	Developing	1644924	45.3	48	160
Antigua and Barbuda	2000	Developing	NA	73.6	156	0
Argentina	2000	Developing	3757452	74.1	137	14
Armenia	2000	Developing	369588	72.0	142	1
Australia	2000	Developed	19153	79.5	78	2

Each row is a **different record**, or observation

2. Manipulate your data using **data verbs**



3. Example

“What countries have the highest life expectancy, on average across time?”

```
(data
%>% select ( country, year, life_expectancy )
%>% filter ( !is.na( life_expectancy ) )
%>% group_by ( country )
%>% summarise ( life_expectancy_average = mean( life_expectancy ) )
%>% arrange ( desc( life_expectancy_average ) )
%>% head ( n = 10 ) )
```

Data verb	Description	In R / dplyr	In SQL	In python (pandas or polars)
SELECT	Selects columns	<code>select()</code>	SELECT	<code>.select()</code>
MUTATE	Adds new columns	<code>mutate()</code>	ADD COLUMN	<code>.mutate()</code>
FILTER	Filters rows based on column values	<code>filter()</code>	WHERE	<code>.filter()</code>
SUMMARISE	Summarises values over rows	<code>summarise()</code>	(aggregate functions)	<code>.agg()</code>
GROUP BY	Group the data frame by column values	<code>group_by()</code>	GROUP BY	<code>.groupby()</code>
SORT or ARRANGE or ORDER	Sort rows by column values	<code>arrange()</code> Or in groups: <code>arrange(.by_group = TRUE)</code>	ORDER BY	<code>.sort()</code>
JOIN	Joins two data frames by shared column values	<code>inner_join()</code> , <code>full_join()</code> , <code>left_join()</code> , <code>right_join()</code>	INNER JOIN etc.	<code>.join()</code>
SLICE	Takes a subset of rows by position	<code>head()</code> and <code>tail()</code> Or in groups: <code>slice_head()</code> , <code>slice_tail()</code> <code>slice_sample()</code>		