

DATA SCIENCE NOTES

Data Wrangling

Verb

Action

arrange()	arrange the rows by column
filter()	filter a subset of rows
Select()	select a subset of columns
mutate()	mutate or create a column
Summarize()	calculate a numerical summary
group-by()	group rows by column

Goals of Wrangling

- Keep certain variables
- define new variables
- reformat/clean variables
- combine datasets
- Process strings/text
- numerically explore datasets

Filtering

Symbol meaning

==	equal to
!=	not equal to
>	greater than
>=	" " and equal to
<	less than
<=	" " and equal to

%in%, c(..., ...) | a list of multiple variables

Dates

mutate to dates from strings

year(today)
month(today) →
week(today)
mday(today)
yday(today)
wday(today) →

TRUE
label = TRUE

Reshaping

Pivoting

Pivot-longer(data, cols,
name_to = "name", values_to = "value")
Combine values into 1 variable

Pivot-wider(data,
names_from = "name",
values_from = "value")

Spread values across new variables

Joining

Mutating Joins

left-join() keeps all observations from the left but discards right ones that don't match
Inner-join() keep only observations that match in both right and left
full-join() keeps all observations

Filtering Joins

Semi-join() discards observations in left without match in right. If multiple matches in left keeps only 1.
anti-join() Discards observations in left that don't match in right.

Factors

order of factor levels + labels

fact_relevel() manually reorder levels
fact_reorder() reorder by another variable
fact_freg() highest to lowest frequency
fact_rev() reverse current order
fact_recode() manually change levels
fact_lump() group least common values

Strings

string functions

str_replace() finds first part and replaces
str_replace_all() replaces all instances
str_to_lower() upper → lower case
str_sub() only keeps a subset of the string
str_length() number of characters
str_detect() if "-" is in the string = True

Data Import: read_csv(), read_delim(), read_sheet(), st_read()