

Process:

End goal is to have data that is ready for analysis (clean so that we don't have "junk in junk out")

① Tidy data; R term but the concept extends to all software

↳ put the data in a format that makes your analysis easier
(ISA 245 \Rightarrow data model)

R \rightarrow

- every row is an observation
- every variable is in one col
- every cell a single value

② Technically correct data

↳ data understood correctly by the software

reasonable
col name

correct classes
for columns

categorical
variable
correct labels

③ consistent data

→ no outliers

→ data that makes sense based on your knowledge of the problem

within a col
e.g. not having
a negative
age

across
cols
in our data
 $\text{count} = \text{registered} + \text{casual}$

most
of the
time

no missing data