

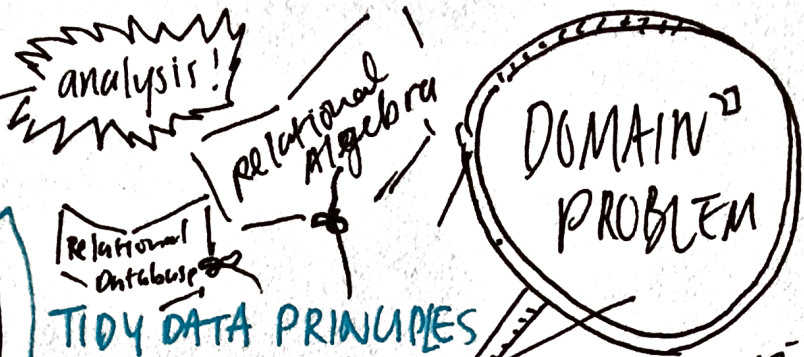
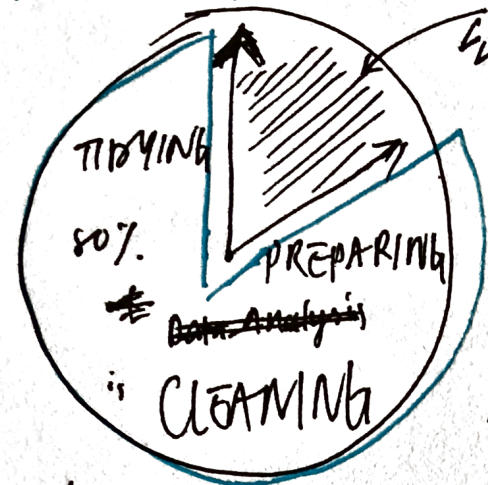
Tidy Data

Hadley Wickham

- # dataCleaning
- # dataTidying
- # relationalDatabases
- # Rstats

Journal of
Statistical
(2014) Software
VOL: 59 NO: 1
PAGES: 1 - 23

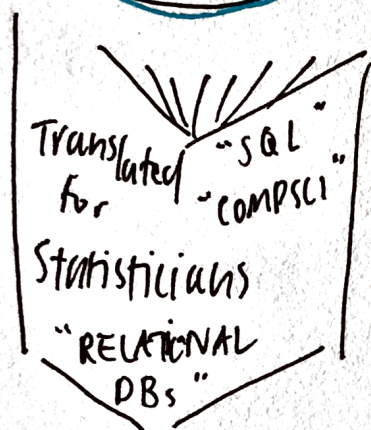
INTRODUCTION.



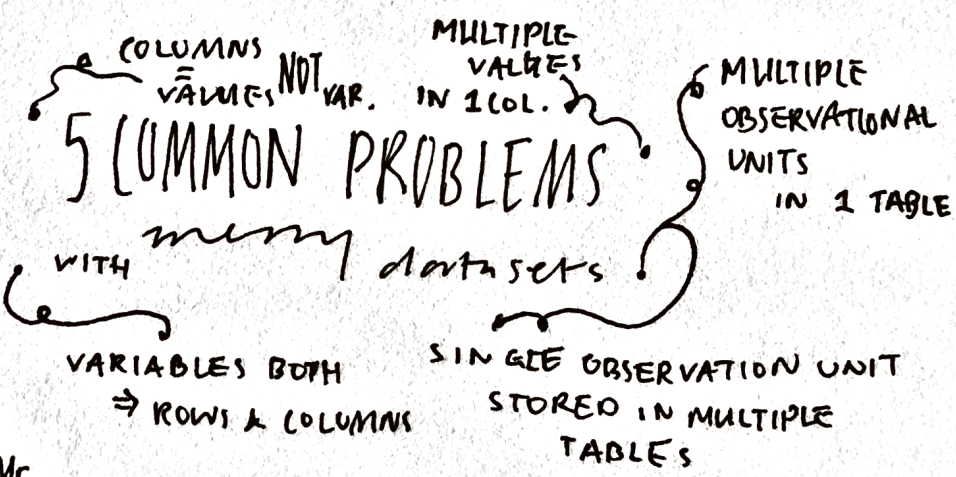
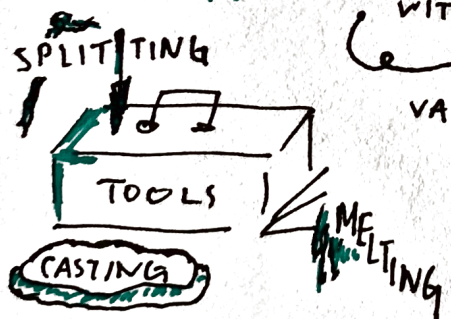
TIDY DATA PRINCIPLES

+ TOOLS

- ☐ facilitate exploration & initial analysis
- ☐ simplify development of data analysis tools that work well together.



TIDYING messy DATA SETS



DEFINING TIDY DATA

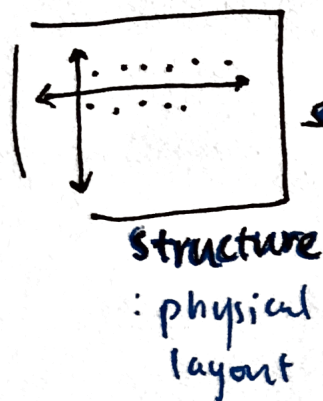
A standard way of mapping meaning to structure.
requires satisfying 3 conditions:

“... every messy dataset is messy
in its own way.”

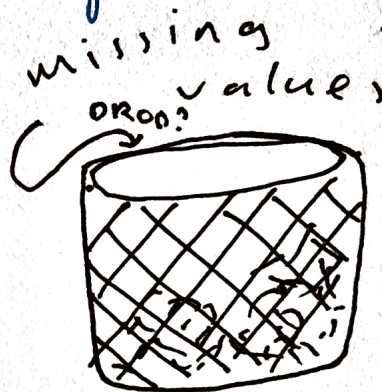
① VARIABLE \Rightarrow COLUMN

② OBSERVATION \Rightarrow ROW

③ UNIT TYPE \Rightarrow TABLE



semantics:
meaning.
var 1 = person
with name
~ first &
~ last...



IS IT AN
OBS. THAT
SHOULD HAVE
BEEN MADE?
YES!
DON'T
DROP
NO
IT'S IMPOSSIBLE
TO MEASURE
 \downarrow
BIN IT!

DATA SET: collection of ⁿVALUES^{strings}

VARIABLE: all values that ^mmeasure the same ^bATTRIBUTE^y
^rUNDERLYING^h

OBSERVATION: all values measure on the
SAME UNIT $\left\{ \begin{array}{l} \text{person} \\ \text{date} \\ \text{location} \end{array} \right.$

ORDER

\rightarrow column \rightarrow
Fixed variables \rightarrow
 $\left\{ \begin{array}{l} \text{MEASURED variables} \\ \text{sort} \\ \text{row by fixed (1, 2, 3...)} \end{array} \right.$