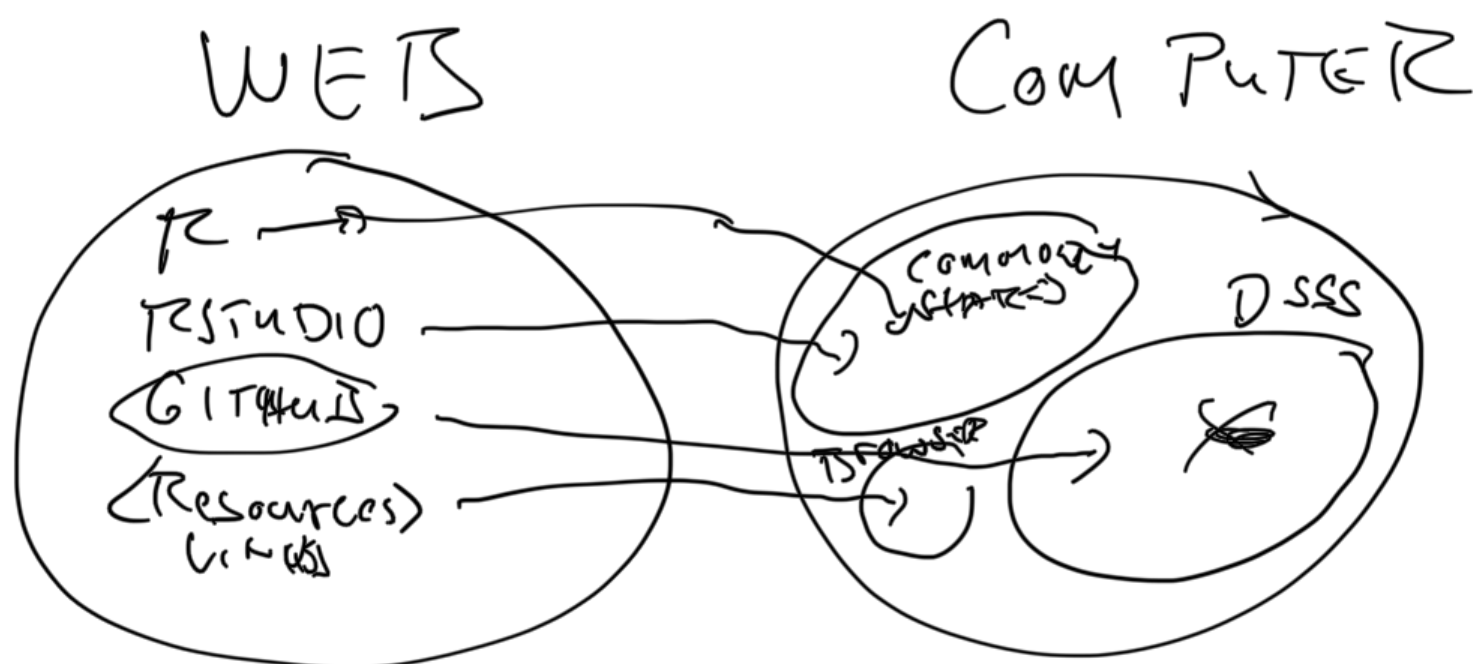


Data science for social sciences



Organization and (File) Names

- (1) DSSS.Rproj
DSSS_01.Rmd
data
literature
figures

(2) (File) Names

- machine readable
- human readable
- play well with default ordering

Henry Bryan
Reproducible
Science Workshop
(see PDF in
literature)

* Machine readable *

- no spaces
- no punctuation (except end)
- no accents or umlaut

a^u, δ^a, u^1

- think about case and be consistent
- deliberate use of delimiters

"_" underscore: delimits unit of metadata that are used later

	1	2
A 1	A1B1	A1B2
2	A2B1	A2B2

Subj: A1_B1 A1_B2 A2_B1 A2_B2
SO1 B.2 2.8 Y.8 J.2

SO1_K2 Sq1-44 SP2-104

"-" hyphen delimit words

2022-02-14

J. Bryan: "so my eyes don't bleed"

* human readable

- name contains info about the content

NOT: my data.csv

data.csv

metadata.csv

analysis ID

Poshew.xlsx

Uliel-01.xlsx

x plays well with default order

- put number first

- use ISO 8601 standard for date

2022-02-14 not: 14.02.2022

- not: S1, S2, ..., S10, ..., S160
But: S001, S002, ..., S010, ..., S100

! WIDE format:

One row of data for the ^{primary} meaningful
unit of analysis

- Subjects

- Country

- Family

**RANDOM
FACTOR**

Child_2020.xlsx
Child School Bridge
2020-02-14

Tsdate T1, ...
2020-02-05

C 00001 SS20
 C 20000 SS20

Sched 2001 Community Pub/Pu
 C 0000

SS20

Community

Questionnaire - VIDE FORM

Subj ^{Item} Item_01 Item_02 ... Item_10 ... Item_33

S01

S02

S60

... WITH STRUCTURE

→ Biography, Selfrating, Opinion
 Subj I01 ... I08 I09 ... I18 I19 ... I33

S01

S02

S60

SAME ... FORM

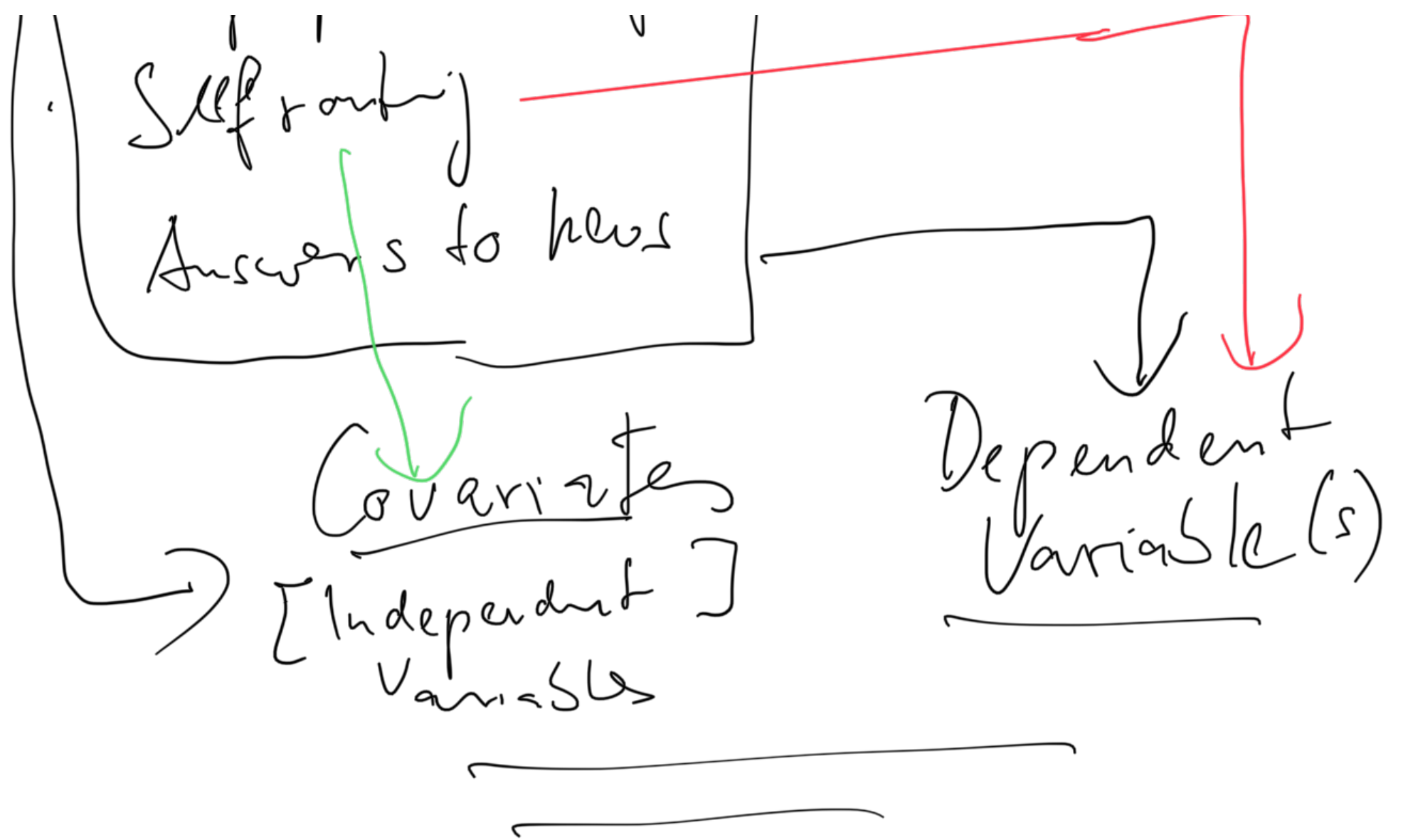
QUESTIONNAIRE

(33 rows / Subj)

Subj	Topic	Item	Response
S01	Bio	I01	male ← Gender
S01	Bio	I02	20 ← Age
S01	Bio	I08	← Educ.
S01	STR	I09	3 ← Life Satisfaction (1-5)
S01	STR	I10	4 ← Health (1-5)
S01	Open	I19	1 ← Credibility (1-6)
S01	Open	I33	6
S02	Bio	I01	female ← Gender
S02	Bio	I02	23 ← Age
S02	Open	I33	5
S03	Bio	I01	male
S60	Open	I33	4

NOTE: This is usually not the optimal long format! Can you see the problem?

if Biographical info



"URDATE"

- cleaned up data
- minimal information, nothing that can be computed from the minimal information
- Variable names
- Format used: CSV, excel, spss, parquet
- Archival for eternity
- un plausible
 - (1) obvious coding error; correct value can be inferred → CORRECT
 - (2) obvious measurement, i.e.

error; correct value cannot be
→ NA (missing) inferred
(?) extreme values
→ leave as is

⇓
PREPROCESSING
(Project_preproc.Rmd)

- Read data
- Read file & describe variables
- Select variables & usually a subset
- Rename variables & see below
- Filter values & remove outliers
- Mutate variables
 - transform variables
 - generate new variables
 - ...
- Relocate variables in columns
- Arrange rows

Iteration

⇓
CHECKS

- Distributions of variables

- boxcox()
- Table of demographic info
 - Table of means, smooths
 - Plots



MULTIVARIATE STATISTICS

- Contrast specification
 - Select model supported by data — not overparameterized
 - Check model residuals — Iteration
 - Fit model
- Iterations must NOT be informed by theoretically relevant statistics (e.g. p-values)



DOCUMENT RESULT

- Tables of coefficients
- Visualize effects
- Move into manuscript or supplement

Next: Variable Names

Next: Factors / Covariates

- ~ Random vs fixed
- ~ Between vs within
- ~ Experimental vs. quasi-experimental