

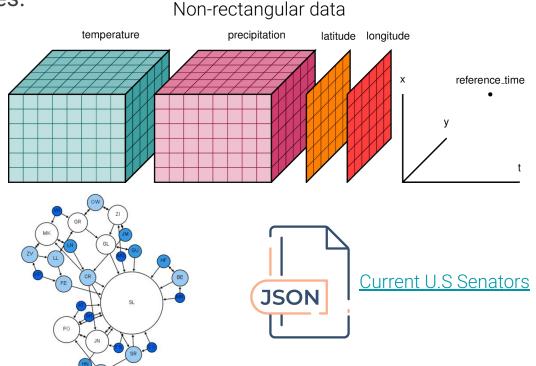
What's in a Dataset?



Rectangular and non-rectangular data

Data come in many different shapes.

Rectangular data

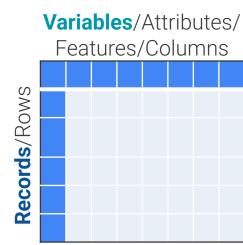




Rectangular data

We prefer rectangular data for analysis (why?)

- Regular structures are easy to manipulate and analyze
- A big part of data cleaning is about transforming data to be more rectangular
- Two kinds of rectangular data: tables and matrices

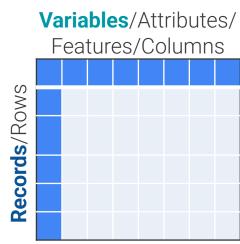




Rectangular data

Tables (a.k.a. dataframes in R/Python and relations in SQL)

- Named columns with different types
- Manipulated using data transformation grammars (e.g., dplyr, pandas)



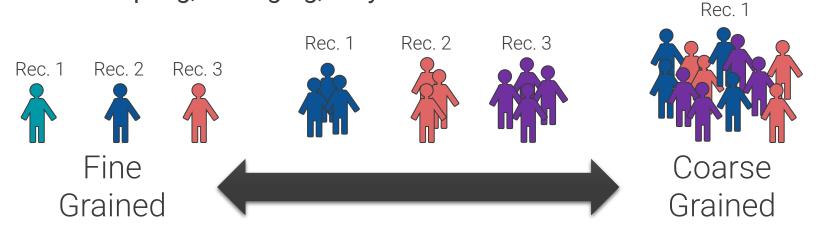
Matrices

- Numeric data of the same type
- Manipulated using linear algebra



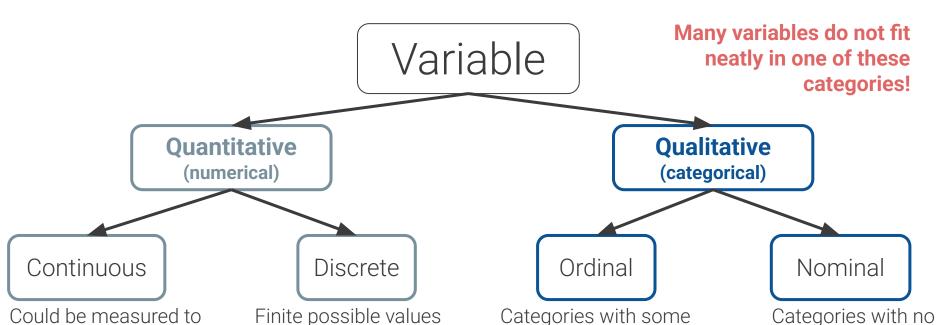
Granularity: how fine/coarse is each record?

- What does each record represent?
 - Examples: a purchase, a person, a group of users
- Do all records capture granularity at the same level?
- If the data are coarse, how were the records aggregated?
 - Sampling, averaging, maybe some of both...





Variables have types



Examples:

- Price
- Temperature

arbitrary precision.

Examples:

- Number of siblings
- Number of classes

Examples:

Preferences

semantic ordering.

Level of education

specific ordering.

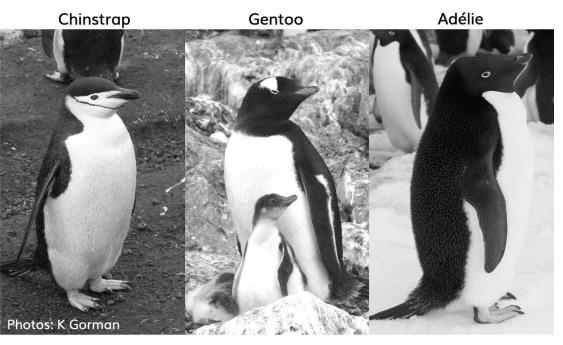
Examples:

- Political Affiliation
- UM C Number



Palmer Penguins







Palmer Penguins



Size measurements for adult foraging penguins near Palmer Station, Antarctica

- Species: Penguin species
- Island: Island in Palmer Archipelago
- Bill length (mm)
- Bill depth (mm)
- Flipper length (mm)
- Body mass (g)
- Sex
- Year: study year

Read more about the data here.



penguins is rectangular data

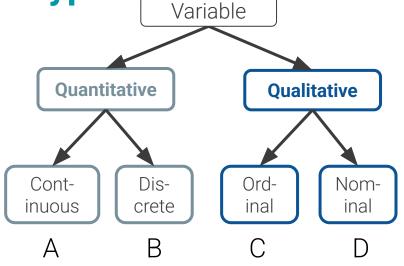
Records represent individual penguins of same granularity.

A tibble: 344	1 × 8						
species <fctr></fctr>	island <fctr></fctr>	bill_length_mm <dbl></dbl>	bill_depth_mm <dbl></dbl>	flipper_length_mm <int></int>	body_mass_g <int></int>	sex <fctr></fctr>	yea
Adelie	Torgersen	39.1	18.7	181	3750	male	200
Adelie	Torgersen	39.5	17.4	186	3800	female	200
Adelie	Torgersen	40.3	18.0	195	3250	female	200
Adelie	Torgersen	NA	NA	NA	NA	NA	200
Adelie	Torgersen	36.7	19.3	193	3450	female	200
Adelie	Torgersen	39.3	20.6	190	3650	male	200
Adelie	Torgersen	38.9	17.8	181	3625	female	200
Adelie	Torgersen	39.2	19.6	195	4675	male	200
Adelie	Torgersen	34.1	18.1	193	3475	NA	200
Adelie	Torgersen	42.0	20.2	190	4250	NA	200

2251603

Review: penguins and variable types

Q	Variable	Data Type
1	species	D. Qualitative Nominal
2	island	D. Qualitative Nominal
3	bill_length_mm	A. Quantitative Cont.
4	body_mass_g	B. Quantitative Discrete
5	sex	D. Qualitative Nominal
6	year	C. Qualitative Ordinal



These examples show how "shaggy" these categories are!
We will revisit variable types in our data visualization unit.