

Data is (still) Messy

Colin Jemmott
DSC 96

Much of this is adapted from the outstanding “Quartz Bad Data Guide”

<https://github.com/Quartz/bad-data-guide>

Last week: identifying messy data

- Are the data types correct?
- String type fields are have consistent values?
- No missing values that we don't understand?
- All values look in a reasonable range?

The data was perfect, right? HA!

How do we deal with the messiness we found?

Last week: identifying messy data

- Are the data types correct?
 - Mostly. Did a little convenience conversion
- String type fields are have consistent values?
 - Case Type, Sex, Ethnicity
 - Solutions: Re-map values (calculated field), filter values, etc...
- All values look in a reasonable range?
 - Age
 - Solutions: filter, smooth,...
- No missing values that we don't understand?
 - Age, Time, Search, Arrested,....
 - Solutions: filter, imputation, create a new binary variable

Human entered data

The dog licensing website for Cook County, Illinois gave a text field to type your dog breed into. As a result this database contained at least 250 spellings of Chihuahua!

How can this be fixed?

One solution: limit choices

SEARCH FOR A BREED

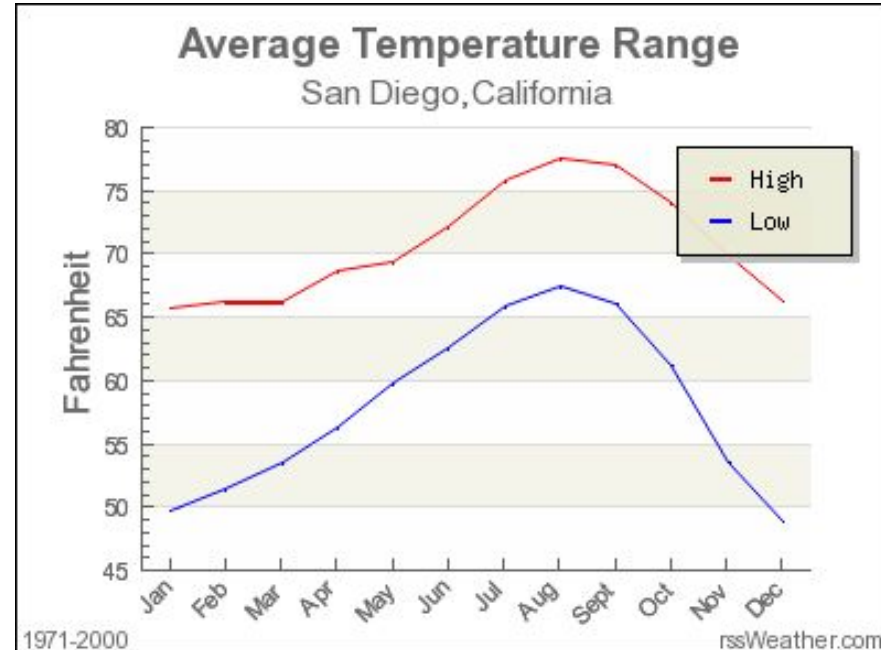
Select A Breed

- Affenpinscher
- Afghan Hound
- Airedale Terrier
- Akita
- Alaskan Malamute
- American English Coonhound
- American Eskimo Dog

Non-Stationary Data

The average low temperature in San Diego is 57 F (14 C). If it is July do you need to bring a sweater?

Sheldon graduated from UCSD CSE in 2010 and got an entry level job paying \$60,000. After working his way up, he is now earning \$68,000. That is more money, right?



Outliers and “Incorrect” Values

- Consistently “nonsense” values
 - Is it a product of the data ingestion process? Time field has year 1899? Is it an inferred “default” value?
 - Solution: Change the value to the correct one!
- Abnormal artifacts from the data collection process
 - E.g. unreasonable spikes in recorded ages at round numbers (25, 35, 45)
 - Solution: Try “smoothing” (e.g. binning the ages)
- Unreasonable outliers
 - Data points with unrealistic and highly unreasonable values. E.g. age=200
 - Solution: filter it? Maybe it points to bugs in the data collection? Maybe it’s **real** and you should investigate!

Understanding How Data is Absent

- Missing by Design (MD)
 - The field being absent is deterministic.
- Missing Completely at Random (MCAR)
 - The missing value isn't associated to the (actual, unreported) value itself, nor the values in any other fields.
 - Example: additional questions in a survey are posed on a random sample of respondents
- Missing at Random (MAR)
 - A missing value may depend only values of other fields, but not its own
 - Example: service workers are less likely to report income.
- Not Missing at Random (NMAR)
 - A missing value depends on the value of the (actual, unreported) variable that's missing.
 - Example: people with high income are less likely to report income.

Null Values: MD, MCAR, MAR, NMAR?

- Attrition due to natural processes?
- Built into the data collection process (intentional)?
- Random issues in (the mechanics of) the data collection process.
- Non-response or refusal

It's very tricky to distinguish between these with certainty!

Can you come up with examples from SDPD dataset?

Null Value Imputation (what to do about them)

- Missing by Design
 - Fill them in? Drop them? Recode the variable?
- Missing Completely at Random (MCAR)
 - Dropping them is ok (if there aren't too many)
- Missing at Random (MAR)
 - Careful! Dropping data will skew your dataset!
 - Replace with mean/mode (perhaps by an associated group)
 - Train a model to replace the missing values
- Not Missing at Random (NMAR)
 - Difficult! Proceed with caution!
 - Train a model to replace the missing values