Lecture 16: Data Context and Quality

CMSC 25910 Spring 2024 The University of Chicago



Data Challenges

Pitfalls of Data Repurposing

- Data is continuously repurposed
 - That's one reason to keep accumulating it
 - ML Training datasets and much more
- Beware of the purpose for which you are repurposing data
 - Why was that dataset created? What was its intended purpose?

Data Quality

- Context
 - Documentation
 - Provenance
 - Assumptions
- Content
 - Errors
 - Missing Data
 - Data formats

Documenting Data

Why we need context

- Data acquisition
- Data stewards
- Data owners
- Data engineers
- Data analysts
- Data consumers

With multiple people involved in the process of transforming raw data into insights, some assumptions with downstream impact may end up buried in the complexity of organizations. Context is important!

Metadata Questions

"We have extracted 155 DGIC questions data workers often face that would be addressed with access to the right MIs. These questions illustrate the metadata landscape we consider in this paper. We have synthesized the 155 questions into 27"

	Representatives of Common Data Questions	DGIC Category	5W1H+R Category
Q1	For what purpose was the dataset created?	G [21, 49]	Why
Q2	Are there tasks for which the dataset should not be used?	G,C [21]	Why
Q3	Who created the dataset?	G,C [21, 26]	Who
Q4	Who was involved in the data creation process?	G,C [21]	Who
Q5	How can the owner/curator/manager of the dataset be contacted?	G [21]	Who
Q6	What are the privacy and legal constraints on the accessibility of the dataset?	C [38]	Who
27	Is there an access control list for the dataset?	G,D [26]	Who
8	What is the reputation of the creator of a dataset?	G [24]	Who
29	What do the instances of the dataset represent?	D,G,I [21]	What
Q10	What is the size of the dataset?	D,G,I [26]	What
Q11	Are there errors in the dataset?	D,G,I [21, 24, 38]	What
Q12	Does the dataset have missing values?	D,G,I [24]	What
213	What is the domain of the values in this dataset?	D,G,I [30]	What
Q14	If the dataset is a sample of a larger dataset, what was the sampling strategy?	G,I [21]	How
Q15	Does the dataset contain personally identifiable information (PII)?	G,C [4, 49]	What
216	What is the quality of the dataset?	G [3, 4, 13, 39]	What
217	Was any preprocessing/cleaning/labeling of the dataset done?	G [21]	How
Q18	Was data collection randomized? Could it be biased in any way?	G [38]	How
Q19	Is there anything about dataset preprocessing/cleaning that could impact future uses?	G [21]	How
Q20	What is the dataset's release date?	D,G,I [30]	When
Q21	Is there an expiration date for this dataset?	D,G [3]	When
Q22	How often will the dataset be updated?	G,I [21]	When
Q23	When was the data last modified?	D,G,I [26]	When
Q24	How easy is it to download and explore this dataset?	D [24]	Where
Q25	What is the format of the dataset, and what type of repository is the dataset located in?	D [38]	Where
Q26	What is the provenance of this dataset?	I [54]	Relationship
Q27	What other datasets exist in this repository that are related to this dataset?	D,G,I [52]	Relationship

Documentation, Context, Semantics

- Provenance/Lineage
 - How was this data recorded/obtained/acquired/produced?
- Metadata
 - Is there documentation associated with the data?
 - What do the attributes mean?
 - What units do they use?
 - (If not) find out that information before using the data. Note the assumptions you had to make

Datasheets for Datasets

Movie Review Polarity

Thumbs Up? Sentiment Classification using Machine Learning Techniques

Motivation

For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

The dataset was created to enable research on predicting sentiment polarity—i.e., given a piece of English text, predict whether it has a positive or negative affect—or stance—toward its topic. The dataset was created intentionally with that task in mind, focusing on movie reviews as a place where affect/sentiment is frequently expressed.¹

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

The dataset was created by Bo Pang and Lillian Lee at Cornell University.

Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number. Funding was provided from five distinct sources: the National Science Foundation, the Department of the Interior, the National Business Center, Cornell University, and the Sloan Foundation.

Any other comments?

None.

Composition

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

The instances are movie reviews extracted from newsgroup postings, together with a sentiment polarity rating for whether the text corresponds to a review with a rating that is either strongly positive (high number of stars) or strongly negative (low number of stars). The sentiment polarity rating is binary {positive, negative}. An example instance is shown in figure 1.

How many instances are there in total (of each type, if appropriate)?

There are 1,400 instances in total in the original (v1.x versions) and 2,000 instances in total in v2.0 (from 2014).

Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample ropresentative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

The dataset is a sample of instances. It is intended to be a random sample of movie reviews from newsgroup postings, with the

¹All information in this datasheet is taken from one of the following five sources; any errors that were introduced are the fault of the authors of the datasheet: http://www.sc.cornell.edu/people/aba/movie-review-data/.http: //xxx.lanl.gov/pdf/cs/0409058v1; http://www.cs.cornell.edu/people/pab/ movie-review-data/t-polaritydata.README 1.0.txt; http://www.cs.cornell. edu/people/pab/movie-review-data/tooldata.README 2.0.txt; these are words that could be used to describe the emotions of john sayles' characters in his latest , limbo . but no , i use them to describe myself after sitting through his latest little exercise in indie egonamia , i can forgive many things . but using some hackneyed , whacked-out , screwed-up * non * ending on a movie is unforgivable . i walked a half-mile in the rain and sat through two hours of typical , plodding sayles melodramat to get cheated by a complete and total copout finale . does sayles think he's roger corman ?

Figure 1. An example "negative polarity" instance, taken from the file neg/cv452_tok-18656.txt.

exception that no more than 40 posts by a single author were included (see "Collection Process" below). No tests were run to determine representativeness.

What data does each instance consist of? "Raw" data (e.g., unprocessed text or images)or features? In either case, please provide a description.

Each instance consists of the text associated with the review, with obvious ratings information removed from that text (some errors were found and later fixed). The text was down-cased and HTML tags were removed. Boilerplate newsgroup header/footer text was removed. Some additional unspecified automatic filtering was done. Each instance also has an associated target value: a positive (+1) or negative (-1) sentiment polarity rating based on the number of stars that that review gave (details on the mapping from number of stars to polarity is given below in "Data Preprocessing").

Is there a label or target associated with each instance? If so, please provide a description.

The label is the positive/negative sentiment polarity rating derived from the star rating, as described above.

Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text. Everything is included. No data is missing.

Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)? If so, please describe how these relationships are made explicit.

None explicitly, though the original newsgroup postings include poster name and email address, so some information (such as threads, replies, or posts by the same author) could be extracted if needed.

Are there recommended data splits (e.g., training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.

The instances come with a "cross-validation tag" to enable replication of cross-validation experiments; results are measured in classification accuracy.

Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description. See preprocessing below.

Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)? If it links

https://www.microsoft.com/en-us/research/publication/datasheets-for-datasets/ (page 1 of 4)

Metadata Management and Catalogs

- Data Catalog:
 - A database for metadata
 - Centralizes tribal knowledge
 - Many challenges to make this work well

Berkeley's Ground [28] Microsoft Azure Data Catalog [31] Apache Atlas [3] Denodo platform [17] SAP Data Intelligence platform [46] Boomi Data platform [7] WeWork's Marquez [50] Lyft's Amundsen [41] Linkedin's Datahub [37]

- Cultural and socio-technical as much as a technical problem
 - Incentives to get people to insert metadata into the catalogs
 - Documenting datasets is not in their critical path except in regulated industries or domains with strong auditors

Data Errors

Types of Data Errors

- Outliers
 - Values that deviate from the distribution (statistical sense)
 - 2, 3, 4, 5654545, 3, 2
- Duplicates
 - Distinct records that refer to the same real-world entity
 - e.g., (first name, last name), (last name, first name)
- Rule Violations
 - Records that violate *integrity constraints*: not null, uniqueness, etc.
- Pattern Violations
 - Violate syntactic and semantic constraints: alignment, misspelling, semantic data types, etc.
 - ZIP code -> State

Types of Data Errors



Figure 2. Classification of data quality problems in data sources

Rahm and Do: Data Cleaning: Problems and Current Approaches. IEEE Data Engineering Bulletin (23): 3-13 (2000)

How Do Errors Affect Analysis?

Dep. Name	Num Employees
Computer Science	41
Economics	112
Statistics	26
CS	41
Physics	33
Chemistry	31

- Duplicates and outliers affect descriptive stats / aggregation
- Outliers are not always errors
 - They may indicate different measurement standards or methods, or particular distributions (e.g., long-tailed)
- Error or not? Can depend on what we are trying to achieve

The Art of Data Cleaning

"As a large mass of raw information, Big Data is not self-explanatory. And yet the specific methodologies for interpreting the data are open to all sorts of philosophical debate. Can the data represent an 'objective truth' or is any interpretation necessarily biased by some subjective filter or the way that data is 'cleaned?" *The Promise and Peril of Big Data. Bollier, 2010, p. 13*

Tooling for Data Cleaning

- OpenRefine
- Ad-hoc tools
- Most data cleaning is accomplished using ad-hoc scripts prepared by data engineers and stewards
 - This is at odds with good documentation of datasets!

- NULL values
 - Many representations: NULL, null, "NULL", "", 0, -1, No, "nil", , NA, Nope...
- NULL values can result from collection or data cleaning
- How do people 'repair' dirty data?
 - A default strategy is to set the value = NULL
- Common approach: Drop rows and hope for the best! 🛞

"...for most of our scientific history, we have approached missing data much like a doctor from the ancient world might use bloodletting to cure disease or amputation to stem infection (e.g, removing the infected parts of one's data by using list-wise or pair-wise deletion). My metaphor should make you feel a bit squeamish, just as you should feel if you deal with missing data using the antediluvian and ill-advised approaches of old." Todd Little. Preface to Applied Missing Data Analysis, Craig Enders.

- Missing data patterns
 - What data is missing (e.g., which cells in a table)
- Missing data <u>mechanisms</u>
 - Aims to find relationships between observed variables and missing data (not necessarily explain why data is missing)

Missing Data Patterns

- Patterns: locations of missing values
- Does not explain why data is missing
- Certain patterns associated with reasons
 - e.g., attrition in multi-phase study



(A) Univariate Pattern







(C) Monotone Pattern



(D) General Pattern



(E) Planned Missing Pattern







Example

- Consider a company's hiring procedure to consist of two stages:
 - IQ test to determine whom to hire
 - Job performance review by a manager 6 months in

	Job performance ratings	
IQ	MAR	
78		
84		
84		
85	—	
87	—	
91	7	
92	9	
94	9	
94	11	
96	7	
99	7	
105	10	
105	11	
106	15	
108	10	
112	10	
113	12	
115	14	
118	16	
134	12	

• Scenario 1

• Why are those values missing?

	Job performance ratings	
IQ	MAR	
78		
84		
84		
85	—	
87	—	
91	7	
92	9	
94	9	
94	11	
96	7	
99	7	
105	10	
105	11	
106	15	
108	10	
112	10	
113	12	
115	14	
118	16	
134	12	

• Scenario 1

- Why are those values missing?
- Missing at Random (MAR). Probability of missing data in attribute X depends on some other attribute Y, but not the values of X.



	Job performance ratings		
IQ	Complete	MCAR	
78	9		
84	13	13	
84	10	_	
85	8	8	
87	7	7	
91	7	7	
92	9	9	
94	9	9	
94	11	11	
96	7	_	
99	7	7	
105	10	10	
105	11	11	
106	15	15	
108	10	10	
112	10		
113	12	12	
115	14	14	
118	16	16	
134	12	_	

Scenario 2

• Why are those values missing?

		Job performance ratin	ıgs
IQ	Complete	MCAR	
78	9	_	
84	13	13	
84	10		
85	8	8	
87	7	7	
91	7	7	
92	9	9	
94	9	9	
94	11	11	
96	7	_	
99	7	7	
105	10	10	
105	11	11	
106	15	15	
108	10	10	
112	10	_	
113	12	12	
115	14	14	
118	16	16	
134	12		

• Scenario 2

• Why are those values missing?

 Missing Completely at Random (MCAR). Probability of missing data in X is unrelated to values of X and unrelated to other attributes.

	Job performance ratings		
IQ	Complete	MNAR	
78	9	9	
84	13	13	
84	10	10	
85	8	_	
87	7		
91	7	_	
92	9	9	
94	9	9	
94	11	11	
96	7	_	
99	7	_	
105	10	10	
105	11	11	
106	15	15	
108	10	10	
112	10	10	
113	12	12	
115	14	14	
118	16	16	
134	12	12	

• Scenario 3

• Why are those values missing?

	Job performance ratings		
IQ	Complete	MNAR	
78	9	9	
84	13	13	
84	10	10	
85	8		
87	7		
91	7	—	
92	9	9	
94	9	9	
94	11	11	
96	7	—	
99	7	—	
105	10	10	
105	11	11	
106	15	15	
108	10	10	
112	10	10	
113	12	12	
115	14	14	
118	16	16	
134	12	12	

• Scenario 3

• Why are those values missing?

• Missing Not at Random (MNAR). Probability of missing data in attribute X is related to the values of X.

	Job performance ratings				
IQ	Complete	MCAR	MAR	MNAR	
78	9			9	
84	13	13	_	13	
84	10			10	
85	8	8	—		
87	7	7			
91	7	7	7		
92	9	9	9	9	
94	9	9	9	9	
94	11	11	11	11	
96	7	—	7		
99	7	7	7		
105	10	10	10	10	
105	11	11	11	11	
106	15	15	15	15	
108	10	10	10	10	
112	10	—	10	10	
113	12	12	12	12	
115	14	14	14	14	
118	16	16	16	16	
134	12	—	12	12	

- Missing at Random (MAR). Probability of missing data in attribute X depends on some other attribute, Y, but not the values of X.
- Missing Completely at Random (MCAR). Probability of missing data in X is unrelated to values of X and unrelated to other attributes.
- Missing Not at Random (MNAR). Probability of missing data in attribute X is related to the values of X.

Handling Missing Data

- Drop rows with any missing data
- Fill in the blanks with the mean (or 0, or a random value)
- Maximum Likelihood Estimation
 - See https://en.wikipedia.org/wiki/Maximum_likelihood_estimation
- Multiple Imputation
 - See https://en.wikipedia.org/wiki/Imputation_(statistics)

Handling Missing Data: Deletion

- Remove tuples that have at least one missing value
 - Assumes MCAR. Otherwise this will bias the data!
 - May reduce sample size a lot!
- Widespread because it's very easy to implement
 - Lots of software packages include a 'drop_null' function
 - See Pandas documentation on 'working with missing data'

Handling Missing Data: Imputation

- Generates a value for each missing data point
- Yields a complete dataset (unlike deletion methods)
- Can produce biased datasets (sometimes even when data is MCAR)
- Arithmetic Mean Imputation/Mean substitution:
 - Reduces the variability of the data -> attenuates standard error/deviation
- Regression Imputation: regression line fit using other (correlated) variable
 - Overestimates correlations
- Stochastic Regression Imputation: Augments regression imputation with a normally distributed residual term (i.e., adds normal noise)
 - Gives unbiased parameter estimates under MAR
- Often requires numerical data; there are advanced techniques for filling categorical data (augmentation, enrichment techniques)

Handling Missing Data: Other Techniques

- Hot-Deck Imputation: Fill missing value with non-missing value
 - Variation: cluster other observations based on variables first
 - Think about the many assumptions this method is making!
- Many other methods:
 - Similar response pattern imputation (similar to hot-deck)
 - Averaging available items
 - Last observation carried forward

Best Practices

- All previous methods assume MCAR and will bias data when data is MAR or MNAR (sometimes even when it's MCAR)
- Maximum Likelihood Estimation (MLE) and multiple imputation can in some cases produce unbiased estimates with MCAR and MAR data, but not with MNAR

Disguised Missing Values

- Phone number:
 - (999)999-9999
- Email address:
 - nope@nope.com
- Age:
 - 666

Source	Table	Column	DMVs
	Diabetes	Blood Pressurse	0
UCI ML	adult	workclass	?
	auun	education	Some College
U.S. FDA	Even Reports	EVENT_DT	20010101, 20030101
data.gov	Vendor Location	Ref_ID	-1
data.gov	Graduation	Regents Num	s, -
data.gov.uk	Accidents 2015	Junction Control	-1

Table 1: Sample DMVs.

From "FAHES: A Disguised Missing Value Detector." KDD 2018

Disguised Missing Values



Disguised Missing Values

"The first time Taylor realized something was amiss was when she received a call in 2011 from a small business owner who angrily blamed her for his customers' email problems...After that initial strange call to Taylor, complaints started pouring in, often with distressing and sometimes criminal accusations aimed at the Arnolds, the Wichita Eagle reported... Officers would show up, accusing them of harboring runaway children. Of keeping girls in the house to make pornographic films. Ambulances appeared, prepared to save suicidal persons. FBI agents, federal marshals and IRS collectors have all appeared on their doorstep."