

Introduction to Data Curation

01.23.2018



- Data
- Curation
- Data Curation
- Data Quality and Data Cleaning
- Metadata and forms of Documentation

Declarative

I know **that** something is the case...

I know that longitude is the geographic coordinate that specifies the east-west position of a point on the Earth's surface.

Often precedes and informs procedural knowledge

I know that XML is an encoding standard for structuring data...

Procedural

I know **how** to do something...

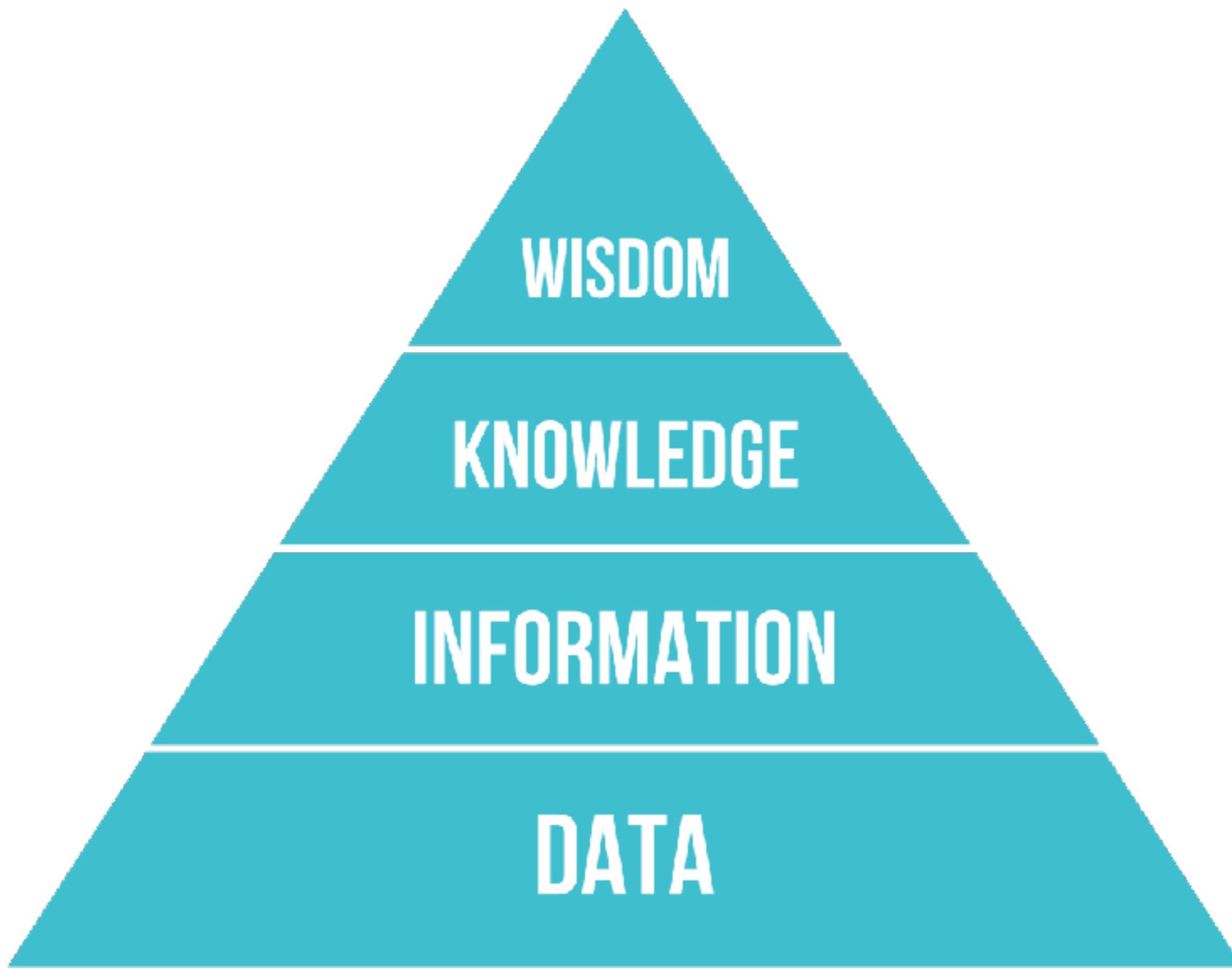
I know how to locate a point on a map given the Longitude and Latitude coordinates...

Often follows from declarative knowledge (but not always)

I know how (and when it is appropriate) to create well-formed XML encodings of free text data.

For this lecture, let's assume that....

Data Curation is the active and ongoing management of **data** throughout a lifecycle of use, including its reuse in unanticipated contexts.



Research Data

“The data, records, files or other evidence, irrespective of their content or form (e.g. in print, digital, physical or other forms), that comprise research observations, findings or outcomes, including primary materials and analysed data.”



Open Data

“Open data is data that can be freely used, shared and built-on by anyone, anywhere, for any purpose.”



For this lecture, let's assume that....

Data are various **types** of information objects playing the
role of evidence.

Evidence of...
Infection (Public Health)
Patterns of Consumer Behavior (Business)
Weather (Atmosphere)

Types vs Roles

Type vs Role distinctions

Type:
Donald Trump is a person.



Role:
Donald Trump is POTUS



Type vs Role distinctions

Type:
XML



Role:
XML as bibliographic data

```
<?xml version="1.0"?>
<catalog>
  - <book id="101">
    <author>Karunaker</author>
    <title>automation Anywhere</title>
    <price>100.20</price>
    <description>book will give info </description>
  </book>
  - <book id="102">
    <author>Rajesh</author>
    <title>Blue prism</title>
    <price>5000</price>
    <description>book will give info </description>
  </book>
  - <book id="103">
    <author>Murali</author>
    <title>Uipath</title>
    <price>100</price>
    <description>book will give info </description>
  </book>
  - <book id="104">
    <author>Balu</author>
    <title>openSpan</title>
    <price>480</price>
    <description>book will give info </description>
  </book>
  - <book id="105">
    <author>Amer</author>
    <title>Workfusion</title>
    <price>10.20</price>
    <description>book will give info </description>
  </book>
</catalog>
```

Data are various **types** of information objects playing the **role** of evidence.

Types of Data

(by file format)



XML



Databases



Flat Files



EDI



Excel



XBRL



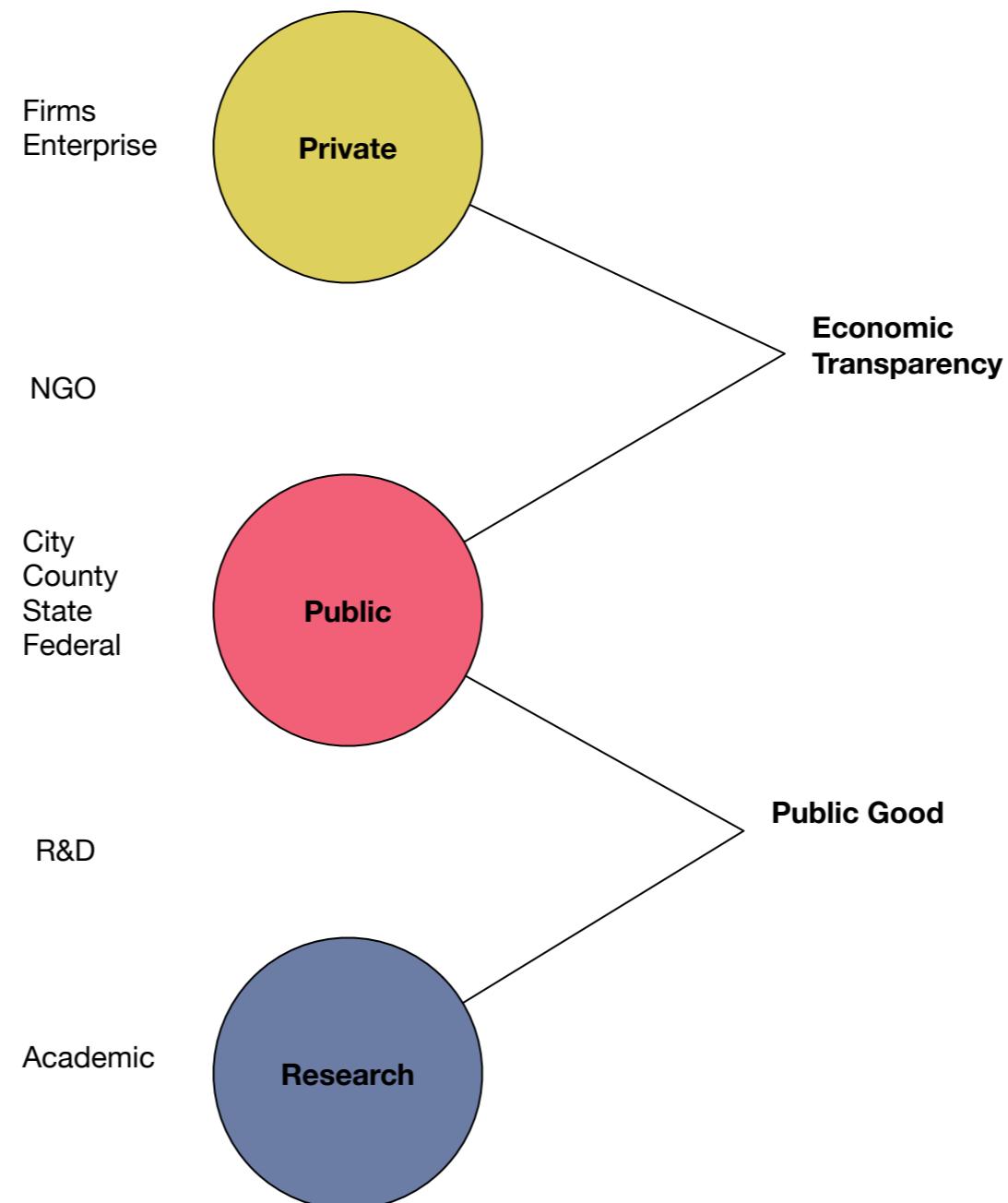
JSON



Web Services

Types of Data

(by sector)



Roles



Roles

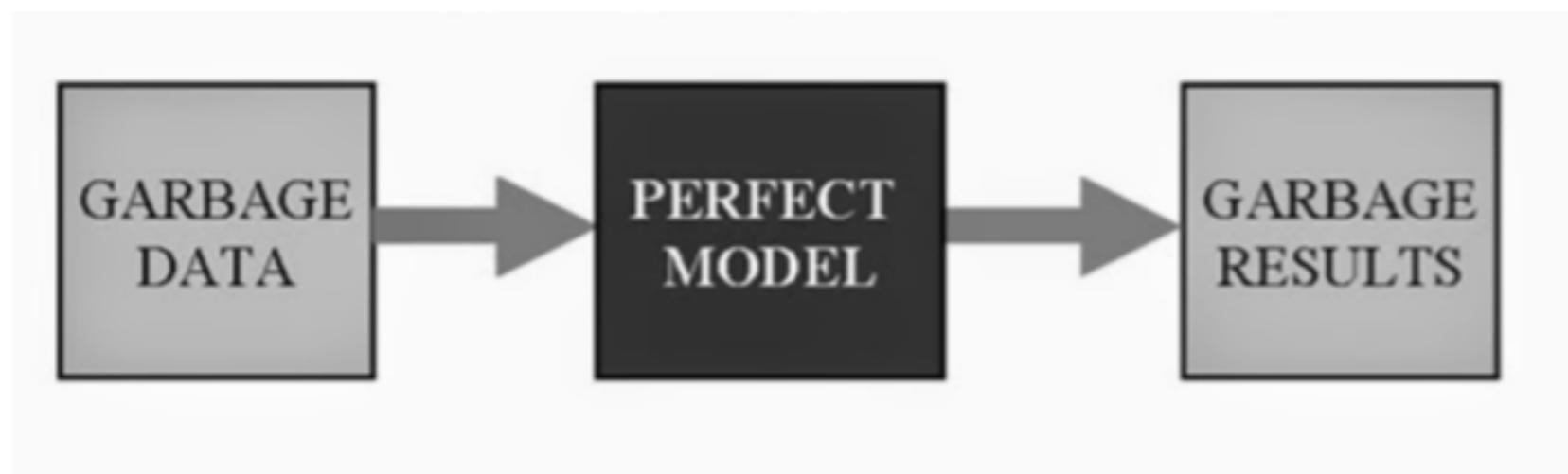


Curation

Data Curation is the active and ongoing management of data throughout **a lifecycle of use**, including its reuse in unanticipated contexts.

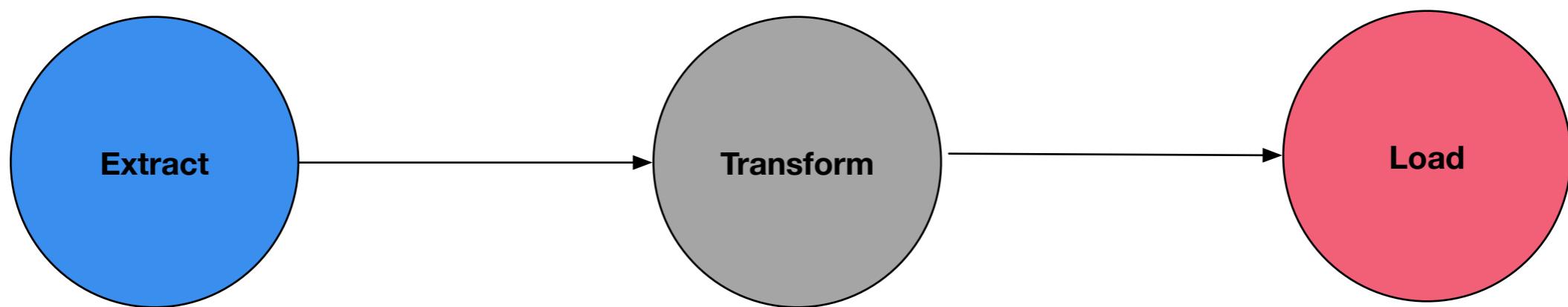
Data Curation

Old computer science saying "Garbage in = Garbage out"

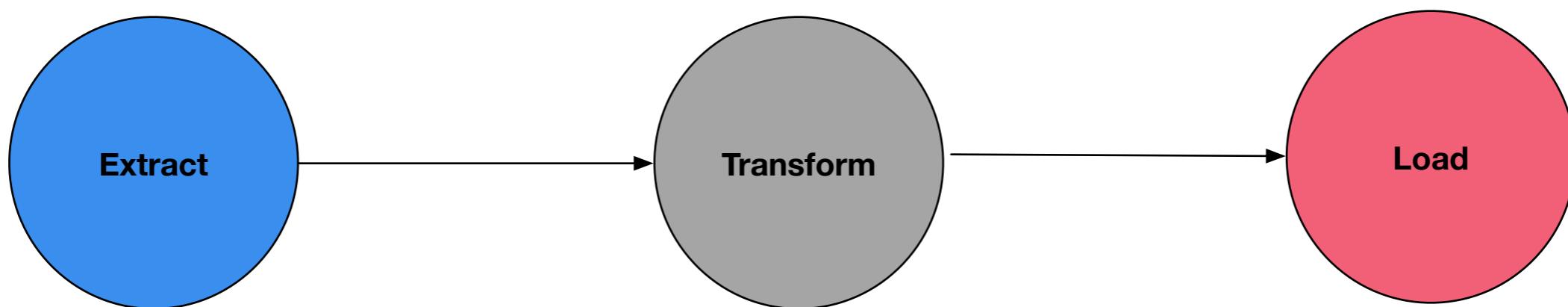


Data curation says "Quality in = Quality out"

ETL: Simplest form of Curation Workflow



ETL: Simplest form of Curation Workflow

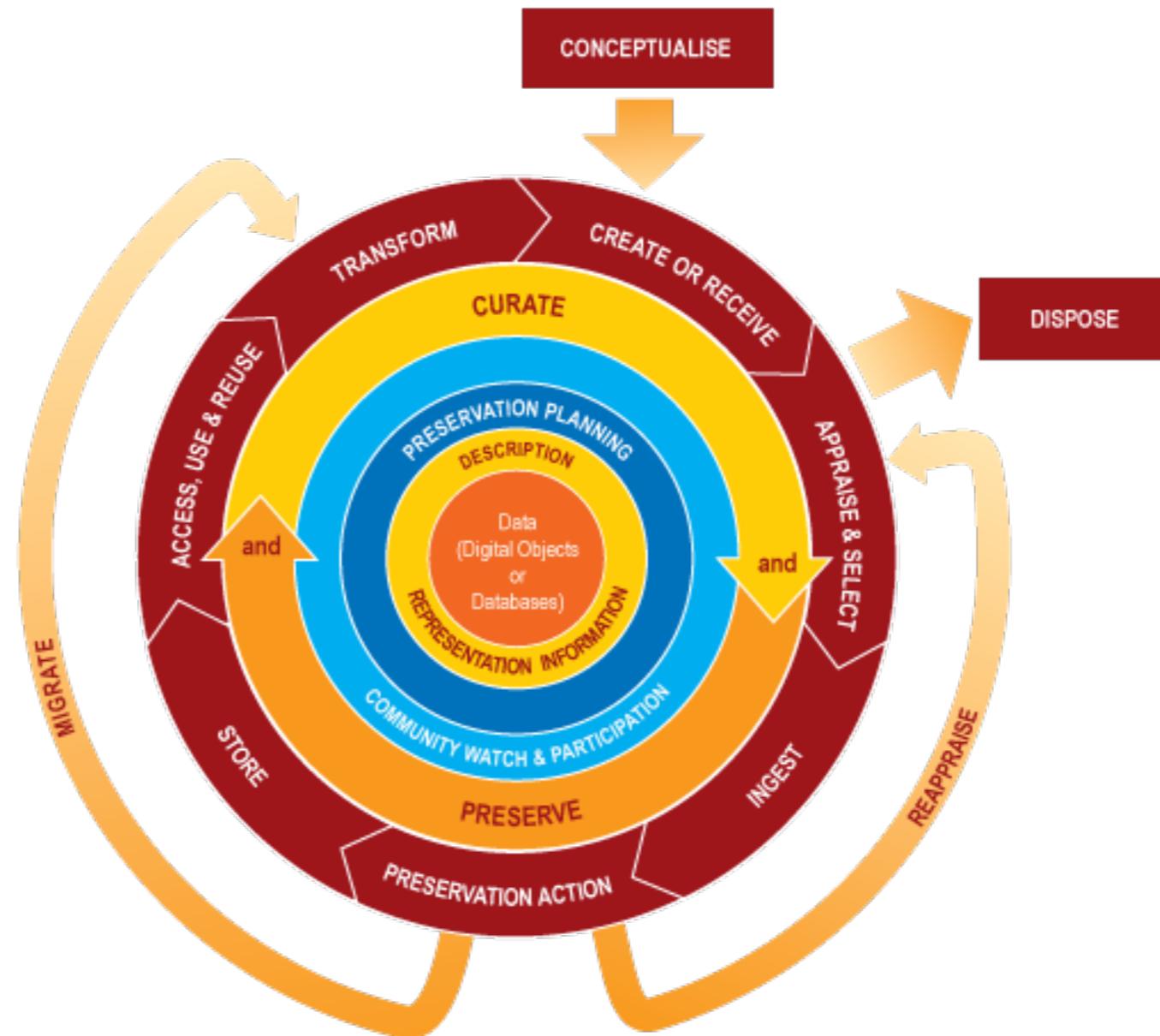


```
<text xmlns="http://www.tei-c.org/ns/1.0" xml:id="d1">
  <body xml:id="d2">
    <div type="book" xml:id="d3">
      <head>Song of Innocence</head>
      <pb n="1"/>
      <div2 type="poem" xml:id="d4">
        <head>Introduction</head>
        <lg type="stanza">
          <l>Piping down the valleys wild,
          <l>Piping songs of pleasant glē,
          <l>On a cloud I saw a child.
          <l>And he laughing said to me
          <l>  Pipe a song about a Lamb.
          <l>So I piped with merry chear.
          <l>Piper pipe that song again.
          <l>So I piped he wept to hear.
          <l>Drop thy pipe thy happy pipe.
          <l>Sing thy songs of happy chear.
          <l>So I sung the same again.
          <l>While he wept with joy to hear.
          <l>Piper sit thee down and write
          <l>In a book that all may read.
          <l>So he vanish'd from my sight.
          <l>And I pluck'd a hollow reed.
          <l>And I made a ratal pen.
          <l>And I stain'd the water clear.
          <l>And I wrote my happy songs.
          <l>Every child may joy to hear.
        </lg>
        <lg type="stanza">
          <l>Piper, sit thee down and write
          <l>In a book, that all may read.
          <l>So he vanish'd from my sight.
          <l>And I pluck'd a hollow reed.
        </lg>
        <lg type="stanza">
          <l>And I made a ratal pen.
          <l>And I stain'd the water clear.
          <l>And I wrote my happy songs.
          <l>Every child may joy to hear.
        </lg>
      </div2>
    </body>
  </div>

```

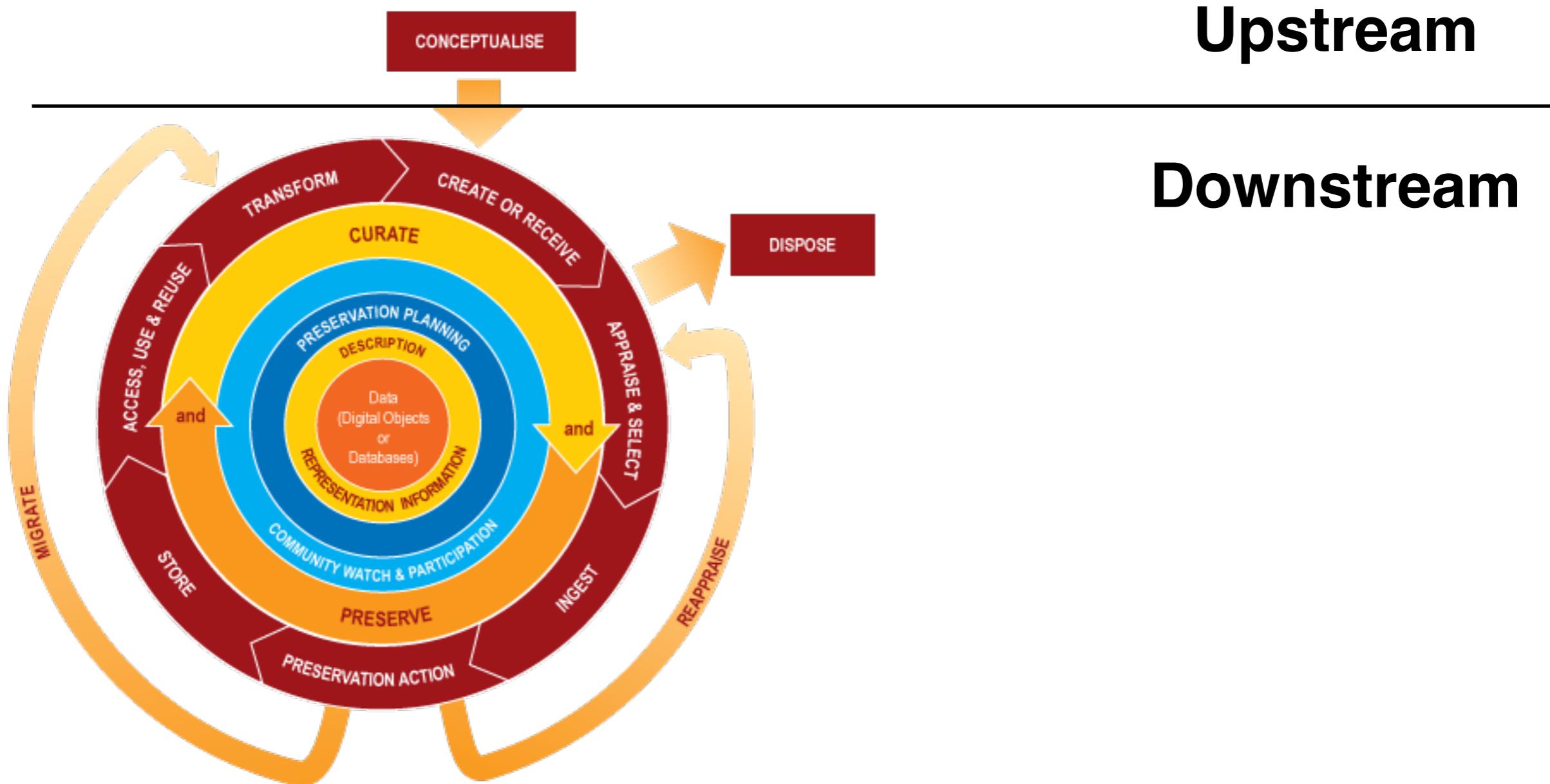


Data Curation Lifecycle



Data Curation Lifecycle

(more complex workflow)



Curation

Data Curation is the active and ongoing management of data throughout a lifecycle of use, including its **reuse in unanticipated contexts**.

Data Quality

*“...the degree to which a set of **characteristics** of data fulfills **stated requirements**.”*

*Examples of **characteristics** are:
completeness, validity, accuracy, consistency, availability and timeliness*

- Normalization
 - Literally - making data conform to a normal schema
 - Figuratively - transforming data structures, organizing variables (columns) and attributes (rows), and editing values so that they are consistent, interpretable, and match best practices in a field.

Database Normalization

(structure)

#	Customer	Order	Item	Delivery Address
1	Linda Porch	01366	Cosa-1	520 Alpha St.
2	Elliott Roof	01377	Ding-1	205 Beta Dr.
3	Kevin Chair	01334	Coisa-1	052 Theta Circle.
4	Todd Window	01355	Veshch-2	502 Gamma Avenue.
5	Diane Door	01353	Koto-2	250 Delta Rd.

A vertical line connects the Customer column of the main table to the Customer column of this secondary table, indicating a one-to-many relationship where multiple delivery addresses are associated with a single customer.

Customer	Delivery Address
Linda Porch	520 Alpha St.
Elliott Roof	205 Beta Dr.
Kevin Chair	052 Theta Circle.
Todd Window	502 Gamma Avenue.
Diane Door	250 Delta Rd.

A vertical line connects the Customer column of the main table to the Customer column of this secondary table, indicating a many-to-many relationship where multiple orders are associated with a single customer.

Customer	Order
Linda Porch	01366
Elliott Roof	01377
Kevin Chair	01334
Todd Window	01355
Diane Door	01353

Data Normalization

(values)

#	Customer	Order	Item	Delivery Address
1	Linda A. Porch	01366	Cosa 1	520 Alpha St.
2	Elliott Roof	1377	Ding-1	205 Beta Drive

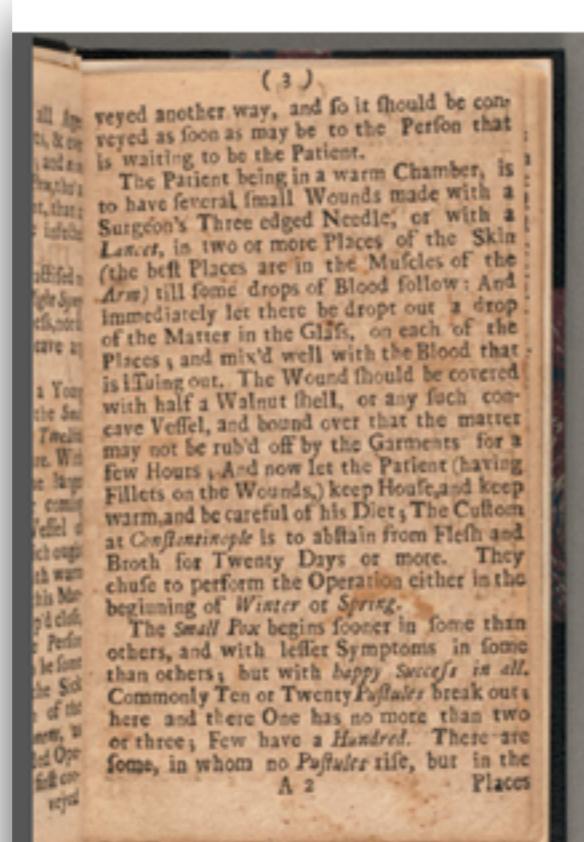
Parsing and standardization — aligning rows and columns, formatting of values into consistent layouts and convention of local standards (for example, postal authority standards for address data). Think of this as the Data Dictionary has to match the Data Value.

Cleaning (aka wrangling / scrubbing) — Modification of data values to meet domain restrictions (e.g. all blank values are to be titled NA), making sure that the **value constraints** of a data dictionary are met in the dataset.

Matching — Identification, linking or merging related entries within or across sets of data (e.g. Item numbers across datasets)

Profiling — Analysis of data to capture statistics that provide insight into the quality of the data and aid in the identification of data quality issues (e.g. We will use a “clustering” feature in Open Refine to group different values together).

Enrichment — Enhancing the value of internally held data. Oftentimes this means adding information (e.g. geographic coordinates; zip codes; etc.)



vveyed another way, and so it should be con-
veyed as soon as may be to the Person that
is waiting to be the Patient.

The Patient being in a warm Chamber, is
to have several small Wounds made with a
Surgeon's Three edged Needle; or with a
Lancet, in two or more Places of the Skin
(the best Places are in the Muscles of the
Arm) till some drops of Blood follow: And
immediately let there be dropt out a drop
of the Matter in the Glass, on each of the
Places ; and mix'd well with the Blood that
is issuing out. The Wound should be covered
with half a Walnut shell, or any such con-
cave Vessel, and bound over that the matter
may not be rub'd off by the Garments for a
few Hours : And now let the Patient (having
Filletts on the Wounds,) keep House, and keep
warm, and be careful of his Diet; The Custom
at Constantinople is to abstain from Flesh and
Broth for Twenty Days or more. They
chuse to perform the Operation either in the
beginning of Winter or Spring.

The Small Pox begins sooner in some than
others, and with lesser Symptoms in some
than others; but with happy Success in all.
Commonly Ten or Twenty Pustules break out;
here and there One has no more than two
or three; Few have a Hundred. There are
some, in whom no Pustules rise, but in the

A 2 Places

File Edit Format View Help

vveyed another way, and fo it shoule be conf-
veyed as foon as may be to the Perfon that
is waiting to be the Patient. _The Patient
being in a warm Chamber, IS to have
i'ev_eral. finall wounds made with a
Sutge'on's Three edged Needle; or" with 2.
Ltmxcr, in two or more Places of the Skin
(the belt Places are in the Mufcles of the
.Ar:w) till foroe drops of Blood follow:
And immediately let there be dropt out a
drop ofthe Matter in the Gldfs, on each 'of
the Places ; and mitfd well with the Blood
that is iffuing out. The Vlfound shoule be
coveted with halfa walnut ihell, or any
fuch con-cave veffel, and bound over that
the matter may not be rub'd off by the
Garments for a few Hours ;~And now let the
Patient (having Fillers on the wounds,)
keep Hou'e,a nd keep warm, and be careful
of his diet ; The cuitotn at
Corjierzfrzoyyle is to abfain from Flefh
and Broth for Twenty days or more. They
chuse to perform the Operation either in
the beginning of wirfrer or spring. ~The
Small Pax begins fooner in forme than
others, and with ie{ller Symptoms in fome
than others; but with ba-ppy Sacfqf: in
all.Commonly Ten or Twentylhjiuffs break
out;here and there one has nomote than two
or three; Few have a Hundred. Thereare
11911161 111 WIOITJ no P#j?efcs ~rife, but
in the A Z ' t ~ _Places t 1 L it h I 1
n

Text Normalization (unstructured data)

- Spelling (e.g. theatre or theater; organise vs organize)
- Vocabulary (e.g.)
- Punctuation (e.g. on-line vs online)
- Chunking (paragraphs, sentences, stanzas, acts, scenes, etc.)
- Markup (what schema did our XML use to encode a text?)



Data Normalization

without writing any code.



Open Refine

FLOSS

Data Cleaning

Looks like a
spreadsheet acts
like a database



R

FLOSS

Statistics + Data
Cleaning

Commandline
power



Microsoft

Proprietary

Analytics + VIZ

Plug + Play



Tableau

Public

VIZ

Plug + Play

PARKING



Can I park here?





R7-1 NO PARKING



R7-2 NO PARKING

What about here...?



[https://data.seattle.gov/Transportation/
SDOT-Street-Signs/atig-uucb#revert](https://data.seattle.gov/Transportation/SDOT-Street-Signs/atig-uucb#revert)

Data Curation is the active and ongoing management of data throughout a lifecycle of use, including its **reuse in unanticipated contexts**.

The reuse of data creates **friction**...
Between person who originally produced the data...
And person trying to understand and use data...

Metadata is a kind of **lubricant** that reduces friction between
data producers and data users



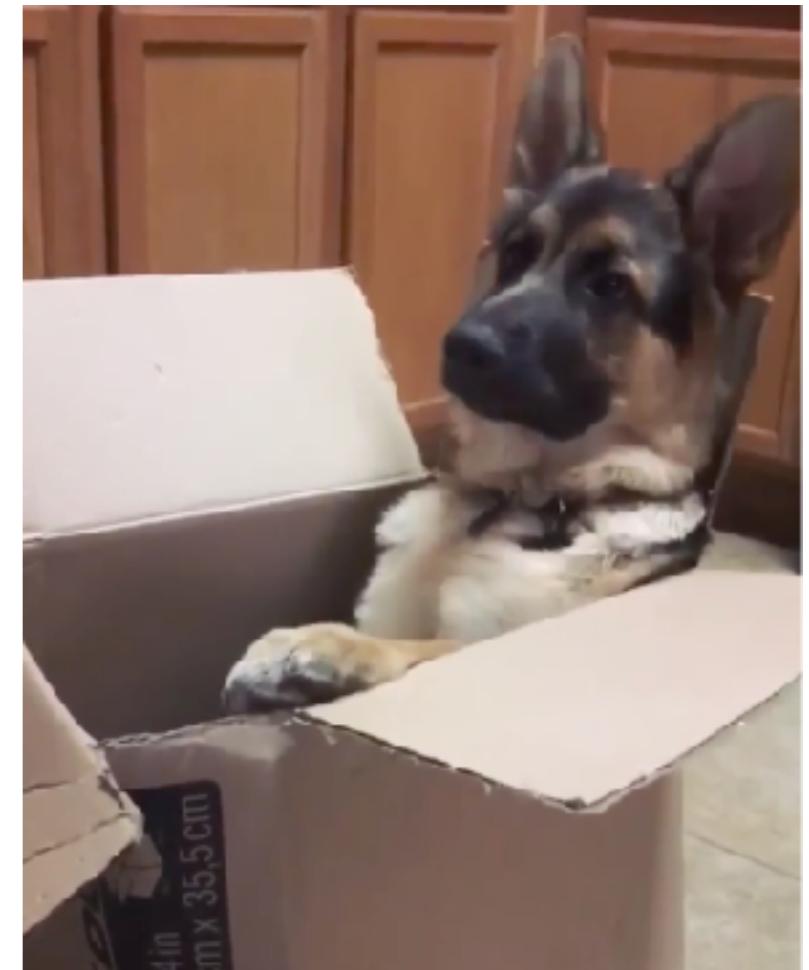
Metadata is most simply a set of **standardized attribute-value** pairs that provide **contextual information** about an object or artifact:

Attribute	Value
Title	Hitchhiker's Guide to the Galaxy
Creator	Douglas Adams

Metadata is most simply a set of **standardized attribute-value** pairs that provide **contextual information** about an object or artifact:

Attributes - properties, features, or characteristics of instances (and by inheritance, classes)

Attribute	Value
Name	Masha
Eye Color	Blue



Expressivity vs. Tractability

(inverse relationship)

The *more expressive* we make our metadata, the *less tractable* it is in terms of generating, managing, and computing for reuse.

The challenge of metadata and documentation for data curation is balancing expressivity and tractability.

Structured vs. Unstructured Metadata

```
1 <metadata xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:dc="http://purl.org/dc/elements/1.1/"*
2   xmlns:dcterms="http://purl.org/dc/terms/" xmlns="http://dublincore.org/documents/dcmi-terms/">
3     <dcterms:title>
4       Replication Data for: Regression Discontinuity with Multiple Running Variables Allowing Partial Effects
5     </dcterms:title>
6     <dcterms:identifier>https://dx.doi.org/10.7910/DVN/PVH6QY</dcterms:identifier>
7     <dcterms:creator>Choi, Jin-Young</dcterms:creator>
8     <dcterms:creator>Lee, Myoung-Jae</dcterms:creator>
9     <dcterms:publisher>Harvard Dataverse</dcterms:publisher>
10    <dcterms:issued>2018-01-14</dcterms:issued>
11    <dcterms:modified>2018-01-14T17:01:39Z</dcterms:modified>
12    <dcterms:description>
13      Data and code to replicate findings in "Regression Discontinuity with Multiple Running Variables Allowing Partial Effects".
14    </dcterms:description>
15    <dcterms:subject>Social Sciences</dcterms:subject>
16    <dcterms:contributor>Choi, Jin-Young</dcterms:contributor>
17    <dcterms:dateSubmitted>2018-01-14</dcterms:dateSubmitted>
18    <dcterms:license>CC0</dcterms:license>
19    <dcterms:rights>CC0 Waiver</dcterms:rights>
</metadata>
```

Machine Readable

The screenshot shows a dataset record titled "Replication Data for Regression Discontinuity with Multiple Running Variables Allowing Partial Effects". The record is identified by doi:10.7910/DVN/PVH6QY. It was published on 2018-01-14. The subject is Social Sciences. The dataset contains data and code to replicate findings in "Regression Discontinuity with Multiple Running Variables Allowing Partial Effects". Contributors are Choi, Jin-Young (Seoul University) and Lee, Myoung-Jae (Korea University). The license is CC0. The rights are CC0 Waiver. The dataset was deposited on 2018-01-14.

Dataset Details	Value
Title	Replication Data for Regression Discontinuity with Multiple Running Variables Allowing Partial Effects
DOI	doi:10.7910/DVN/PVH6QY
Published Date	2018-01-14
Subject	Social Sciences
Description	Data and code to replicate findings in "Regression Discontinuity with Multiple Running Variables Allowing Partial Effects" (2018-01-14)
Contributor	Choi, Jin-Young (Seoul University); Lee, Myoung-Jae (Korea University)
License	CC0
Rights	CC0 Waiver
Deposit Date	2018-01-14

Human Readable

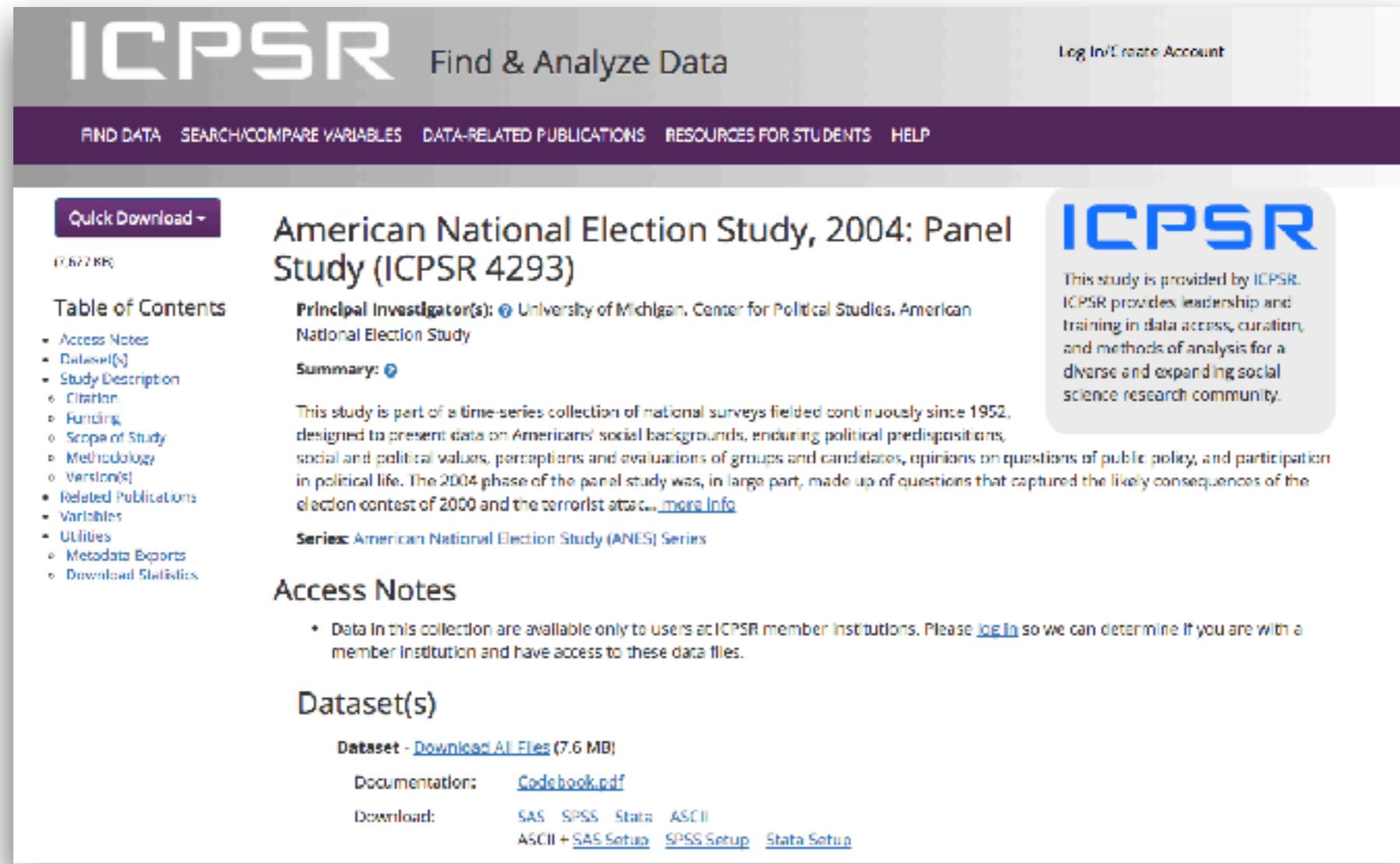
Three basic forms of structured metadata in data curation

Descriptive Metadata: Tells us about objects, their creation, and the context in which they were created (Title, Author, Date)

Technical Metadata: Tells us about the context of the data collection (Instrument, Computer, Algorithm)

Administrative Metadata: Tell us about the management of that data (Rights statements, Provenance, etc.)

Descriptive Metadata: Tells us about objects, their creation, and the context in which they were created (Title, Author, Date)



The screenshot shows the ICPSR dataset page for the American National Election Study, 2004: Panel Study (ICPSR 4293). The page has a header with the ICPSR logo and "Find & Analyze Data". It includes a navigation bar with links for FIND DATA, SEARCH/COMPARE VARIABLES, DATA-RELATED PUBLICATIONS, RESOURCES FOR STUDENTS, and HELP. A "Log In/Create Account" link is also present.

Dataset Summary:

- Quick Download:** 0.672 KB
- Title:** American National Election Study, 2004: Panel Study (ICPSR 4293)
- Principal Investigator(s):** University of Michigan, Center for Political Studies, American National Election Study
- Summary:** This study is part of a time-series collection of national surveys fielded continuously since 1952, designed to present data on Americans' social backgrounds, enduring political predispositions, social and political values, perceptions and evaluations of groups and candidates, opinions on questions of public policy, and participation in political life. The 2004 phase of the panel study was, in large part, made up of questions that captured the likely consequences of the election contest of 2000 and the terrorist attack... [more info](#)
- Series:** American National Election Study (ANES) Series

Access Notes:

- Data in this collection are available only to users at ICPSR member institutions. Please [log in](#) so we can determine if you are with a member institution and have access to these data files.

Dataset(s):

- Dataset:** [Download All Files \(7.6 MB\)](#)
- Documentation:** [Codebook.pdf](#)
- Download:** [SAS](#) [SPSS](#) [Stata](#) [ASCII](#)
[ASCII + SAS Setup](#) [SPSS Setup](#) [Stata Setup](#)

Descriptive Metadata: Tells us about objects, their creation, and the context in which they were created (Title, Author, Date)

```
<codeBook version="2.1" ID="ICPSR04245">
  <docDscr>
    <citation>
      <titlStmt>
        <titl>Metadata record for ANES 2004 Time Series Study</titl>
        <IDNo agency="ICPSR">4245</IDNo>
      </titlStmt>
      <prodStmt>
        <producer abbr="ICPSR">
          <ExtLink URI="http://www.icpsr.umich.edu/images/icpsr-logo.gif" title="ICPSR Logo" role="image"/>
          Inter-university Consortium for Political and Social Research
          <ExtLink URI="http://www.icpsr.umich.edu/ICPSR/" title="URL of ICPSR Web Site"/>
        </producer>
      </prodStmt>
    </citation>
  </docDscr>
</codeBook>
```

<titl> Title

- Mandatory
- Not Repeatable
- Attributes: [ID](#), [xml:lang](#), [source](#)

Description: Full authoritative title for the work at the appropriate level: marked-up document; marked-up document source; study; other material(s) related to study description; other material(s) related to study. The study title will in most cases be identical to the title for the marked-up document. A full title should indicate the geographic scope of the data collection as well as the time period covered. Title of data collection (2.1.1.1) maps to Dublin Core Title element. This element is required in the Study Description citation.

Example(s):

```
<titl>Domestic Violence Experience in Omaha, Nebraska, 1986-1987</titl>
<titl>Census of Population, 1950 [United States]: Public Use Microdata Sample</titl>
<titl>Monitoring the Future: A Continuing Study of American Youth, 1995</titl>
```

Technical Metadata: Tells us about the context of the data collection (Instrument, Computer, Algorithm)

The screenshot shows the PRONOM website interface. At the top, there's a navigation bar with links like 'About us', 'Education', 'Research', 'Information management', and 'Archived sector'. Below that is a search bar. The main content area has a header 'The technical registry PRONOM' with a logo. It displays a 'File format summary' for 'JPEG File Interchange Format 1.01'. The summary table includes fields such as Name (JPEG File Interchange Format), Version (1.01), Other names (JFIF (1.01)), Identifiers (PUID: Anvl<0>, IANA: image/jpeg, Apple Uniform Type Identifier: public.jpeg), Family (Image (Raster)), Classification (Image (Raster)), Disclosure (Full), Description (A detailed description of the format), Orientation (Binary), Byte order (Big-endian (Motorola)), and Related file formats (Raw JPEG Stream, JPEG File Interchange Format (1.02), and JPEG File Interchange Format (1.00)). There are also tabs for Summary, Documentation, Signatures, Compression, Character Encoding, Images, and References.

```
<?xml version="1.0" encoding="utf-8"?>
<PRONOM-Report xmlns="http://pronom.nationalarchives.gov.uk">
  <report_format_detail>
    <FileFormat>
      <FormatID>668</FormatID>
      <FormatName>JPEG File Interchange Format</FormatName>
      <FormatVersion>1.01</FormatVersion>
      <FormatAliases>JFIF (1.01)</FormatAliases>
      <FormatFamilies>
        </FormatFamilies>
      <FormatTypes>Image (Raster)</FormatTypes>
      <FormatDisclosure>Full</FormatDisclosure>
      <FormatDescription>The JPEG File Interchange Format (JFIF) is a file format for storing JPEG-compressed raster images. It was developed by the Independent JPEG Group and C-Cube Microsystems. In the absence of any such format being defined in the JPEG standard, and mainly because a lot of software that is not compliant with the standard is written to be able to read JFIF files, it has become a de facto standard; this is what is commonly referred to as the 'JPEG file format'. A JFIF file comprises a JPEG data stream together with a JFIF marker. It begins with a Start of Image (SOI) marker, immediately followed by a JFIF Application (APP0). This is followed by the JPEG image data, which is terminated by an End of Image (EOI) marker. JFIF supports up to 24-bit colour and uses lossy compression (based on the Discrete Cosine Transform algorithm). Other types of compression are available through JPEG extensions, including progressive image buildup, arithmetic encoding, variable quantization, selective refinement, image tiling, and lossless compression, but these may not be supported by all JFIF readers and writers.</FormatDescription>
      <BinaryFileFormat>Binary</BinaryFileFormat>
      <ByteOrders>Big-endian (Motorola)</ByteOrders>
      <ReleaseDate>
```

Administrative Metadata: Tell us about the management of that data (Rights statements, Provenance, etc.)



```
<premis:object>
  <!--other metadata-->
  <premis:signatureInformation>
    <premis:signatureInformationEncoding>BASE 64</premis:signatureInformationEncoding>
    <premis:signer>Susan Thomas</premis:signer>
    <premis:signatureMethod>DSA-SHA1</premis:signatureMethod>
    <premis:signatureValue>qUADDMHZkyebvRdLs+6Dv7RvgMLRIUaDB4Q9yn9XoJA79a2882ffTg==
    </premis:signatureValue>
    <premis:signatureValidationRules>Add reference to repository documentation detailing signature validation rules</premis:signatureValidationRules>
    <premis:signatureProperties>2006-11-01T10:15:16</premis:signatureProperties>
  </premis:signatureInformation>
  <!--other metadata-->
</premis:object>
```

Unstructured Metadata or Documentation

(human readable)

README.txt - provides narrative explanation of what a dataset contains, how it was produced, and how it can or should be used.

Data Dictionary - defines the variables (and constraints on the values of those variables) in a dataset

CodeBook - defines what codes were created to analyze, or summarize a dataset

readMe.txt

The screenshot shows a Harvard Dataverse dataset page. At the top, there are tabs for Home, About, User Guide, Report, Share My, and Log In. Below the tabs, there are sections for 'About' (1 points), 'Data' (1 datasets), 'Code' (0), and 'Documents' (1). A red box highlights the 'readMe.txt' document. The document preview shows its content, which includes a table of contents and some text. At the bottom, there is a 'File Metadata' table with columns for Name, Description, and Last Modified.

Name	Description	Last Modified
readMe.txt	1.0 MB (1 page) Last Modified: Jan 14, 2018. Information about the replication data for Regression Discontinuity with Multiple Running Variables Allow for Partial Effects. This file contains the Stata do file and the log file.	Jan 14, 2018
DATA	1.0 MB (1 file) Last Modified: Jan 14, 2018. Data files for the regression discontinuity analysis. This file contains the Stata dta file.	Jan 14, 2018
RESULTS	1.0 MB (1 file) Last Modified: Jan 14, 2018. Results from the Stata do file. This file contains the Stata log file.	Jan 14, 2018
CODE	1.0 MB (1 file) Last Modified: Jan 14, 2018. Gauss code for the regression discontinuity analysis. This file contains the Gauss gau file.	Jan 14, 2018

There are two folders to replicate the empirical results of the paper: STATA folder and GAUSS folder.

The STATA folder provides the graphic outputs in *.gph files, and the GAUSS folder provides the table outputs in *.txt files.

Even if the user is unfamiliar with GAUSS, he/she can still obtain at least parts of the table outputs by running the STATA program; specifically, the estimates of the tables in the paper, and the t-values computed with the usual OLS asymptotic variance estimator, but not the confidence intervals (CI's) computed with bootstrap in the paper.

The details of the STATA and GAUSS folders are as follows.

----- STATA FOLDER DESCRIPTION -----

The enclosed STATA program "Election_26AUG2017_Stata.do" produces Table 1, all estimates in Tables 2 and 3, and Figures 2 and 3.

The *.log file is the saved result corresponding to the .do file and it includes Tables 1, 2, and 3.

And the *.gph files are figure outputs also generated with the .do file.

What the STATA program does not produce is the confidence intervals (CI) based on bootstrap in Tables 2 and 3; instead of the CI's, the STATA program provides the usual t-values based on the OLS asymptotic variance estimator for all OLS-based estimates. Because of this, the OLS CI's in the paper differ somewhat from those in the STATA output file.

The STATA program does not provide any t-value for the "boundary-weighting (BW)" estimator in Tables 2 and 3, because BW is a complicated estimator, not based on OLS.

If the reader desires to generate bootstrap CI's, he/she may use the bootstrap option for OLS provided by STATA.

In the STATA program, "mf" appears, which stands for "multiplicative factor" in selecting the bandwidth

$$h = mf \cdot SD(S) \cdot N^{-1/6} \quad \text{where } S \text{ is the running variable in use.}$$

The "mf" value is typically about 3.5–2.5, and it was already chosen with Cross-Validation (CV) using a GAUSS program.

The STATA file uses the pre-selected value of "mf" without redoing the CV procedure.

The reason for not providing the bootstrap CI's and not doing the CV procedure in the STATA program is that these procedures require a sophisticated programming with STATA, which the authors could not do, as they are not regular users of STATA.

----- GAUSS FOLDER DESCRIPTION -----

In the GAUSS folder, all files are written in GAUSS, which is a programming language from Aptech Systems Inc.

GAUSS files can be opened with any text file editor (e.g., notepad or wordpad).

In our paper, empirical parts were done with GAUSS, except for Figures 2 and 3.

Data Dictionary

Department	Dataset Name	Field Name	Field Alias	Field Type	API Key	Field Definition	Field Type Flag
Rent Arbitration Board	Eviction Notices	City	city	text	city	The city where the eviction notice was issued. In this dataset, always San Francisco.	
Rent Arbitration Board	Eviction Notices	State	state	text	state	The state where the eviction notice was issued. In this dataset, always CA.	
Rent Arbitration Board	Eviction Notices	Eviction Notice Source Zipcode	zip	text	zip	The zip code where the eviction notice was issued.	
Rent Arbitration Board	Eviction Notices	File Date	file_date	timestamp		The date on which the eviction notice was filed with the Rent Board of Arbitration.	
Rent Arbitration Board	Eviction Notices	Non Payment	non_payment	boolean	non_payment	This field is checked (true) if the landlord indicated non-payment of rent as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Breach	breach	boolean	breach	This field is checked (true) if the landlord indicated breach of lease as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Nuisance	nuisance	boolean	nuisance	This field is checked (true) if the landlord indicated nuisance as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Illegal Use	illegal_use	boolean	illegal_use	This field is checked (true) if the landlord indicated an illegal use of the rental unit as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Failure to Sign Renewal	failure_to_sign_renewal	boolean	failure_to_sign_renewal	This field is checked (true) if the landlord indicated failure to sign over renewal as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Access Denial	access_denial	boolean	access_denial	This field is checked (true) if the landlord indicated unlawful denial of access to unit as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Unapproved Subtenant	unapproved_subtenant	boolean	unapproved_subtenant	This field is checked (true) if the landlord indicated the tenant had an unapproved subtenant as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Owner Move In	owner_move_in	boolean	owner_move_in	This field is checked (true) if the landlord indicated an owner move in as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Demolition	demolition	boolean	demolition	This field is checked (true) if the landlord indicated demolition of property as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Capital Improvement	capital_improvement	boolean	capital_improvement	This field is checked (true) if the landlord indicated a capital improvement as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Substantial Rehab	substantial_rehab	boolean	substantial_rehab	This field is checked (true) if the landlord indicated substantial rehabilitation as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Ellis Act Withdrawal	ellis_act_withdrawal	boolean	ellis_act_withdrawal	This field is checked (true) if the landlord indicated an Ellis Act withdrawal (going out of business) as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Condo Conversion	condo_conversion	boolean	condo_conversion	This field is checked (true) if the landlord indicated a condo conversion as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Roommate Same Unit	roommate_same_unit	boolean	roommate_same_unit	This field is checked (true) if the landlord indicated if they were evicting a roommate in their unit as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Other Cause	other_cause	boolean	other_cause	This field is checked (true) if some other cause not covered by the acmfa code was indicated by the landlord. These are not enforceable grounds.	
Rent Arbitration Board	Eviction Notices	Late Payments	late_payments	boolean	late_payments	This field is checked (true) if the landlord indicated habitual late payment of rent as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Lead Remediation	lead_remediation	boolean	lead_remediation	This field is checked (true) if the landlord indicated lead remediation as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Development	development	boolean	development	This field is checked (true) if the landlord indicated a development agreement as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Good Samaritan Ends	good_samaritan_ends	boolean	good_samaritan_ends	This field is checked (true) if the landlord indicated the period of good samaritan laws coming to an end as a grounds for eviction.	
Rent Arbitration Board	Eviction Notices	Constraints Date	constraints_date	timestamp	constraints_date	In the case of certain just cause evictions like Ellis and Owner Move In, constraints are placed on the property and recorded by the City.	
Rent Arbitration Board	Eviction Notices	Supervisor District	supervisor_district	numeric	supervisor_district	There are 11 members of the Board of Supervisors in San Francisco, each representing a geographic district. These are numbered 1 through 11.	
Rent Arbitration Board	Eviction Notices	Neighborhoods - Analysis Boundaries	neighborhood	text	neighborhood	The Department of Public Health and the Mayor's Office of Housing and Community Development, with support from the Planning Dept.	
Rent Arbitration Board	Eviction Notices	Location	client_location	Geometry	geometry: p:client_location	Contains the geometry of the record in Well Known Text (WKT) format.	

<https://data.sfgov.org/City-Management-and-Ethics/-alpha-Master-data-dictionary/wn8x-uk7i#>

Codebook

CODEBOOK FOR ICPSR 9028			V5	DIVISION	13	13	F1
UNIFORM CRIME REPORTING PROGRAM DATA (UNITED STATES)			V6	YEAR	14	17	F4
PART 1: OFFENSES KNOWN AND CLEARANCES BY ARREST, 1980			V7	CITY SEQUENCE NUMBER	18	22	F5
PLEASE NOTE: The "M" between the code and the code label indicates the code has been designated as a missing value.			V8	CORE CITY INDICATION	23	23	A1
NAME	VARIABLE LABEL		BEG COL	END COL	FMT		
V1	ID CODE		1	1	F1		
	1 Offenses known						
V2	NUMERIC STATE CODE		2	3	F2	N	No, not core city of MSA
	1 Alabama					Y	Yes, core city of MSA
	2 Arizona						
	3 Arkansas						
	4 California						
	5 Colorado						
	6 Connecticut						
	7 Delaware						
	8 District of Columbia						
	9 Florida						
	10 Georgia						
	11 Idaho						
	12 Illinois						
	13 Indiana						
	14 Iowa						
	15 Kansas						
	16 Kentucky						
	17 Louisiana						
			V9	COVERED BY CODE	24	30	A7
			V10	LAST UPDATE	31	38	F8
			V11	FIELD OFFICE	39	42	F4
			V12	NUMBER OF MONTHS REPORTED	43	44	F2
				0 No months reported			
				1 Jan last reported			
				2 Feb last reported			
				3 March last reported			
				4 April last reported			
				5 May last reported			
				6 June last reported			

- Metadata helps reduce friction between data producers and data users
- Comes in two forms: Structured and Unstructured
- Structured metadata uses an encoding, and a formally defined schema to make metadata **Machine Readable**
- Unstructured Metadata is meant to provide contextual information that is **Human Readable**

Data Curation is the active and ongoing management of data throughout a lifecycle of use, including its reuse in unanticipated contexts.

Data are various **types** of information objects playing the **role** of evidence.

“...the degree to which a set of characteristics of data fulfills stated requirements.”