# Metadata & Documentation

for
Data Curation
**Lecture 2**

INFX 551
Winter 2018

**Lecture 1**

- Ontology & Epistemology

- Ontology (in Information Science)

- Knowledge Representation: Classes vs Instances

    - Attribute Value Pairs

- Expressivity vs Tractability


**Lecture 2**

- Metadata & Documentation

        - Structured vs Unstructured

- Encodings

- Standards Schema Development (application profiles)

- Forms of metadata

- Forms of documentation

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

**<dc: title>**Hitchhikers Guide to the Galaxy**</dc:title>**

| Term Name: | title |
|---|---|
| **URI:** | http://purl.org/dc/terms/title |
| **Label:** | Title |
| **Definition:** | A name given to the resource. |
| **Type of Term:** | Property |
| **Refines:** | http://purl.org/dc/elements/1.1/title |
| **Version:** | http://dublincore.org/usage/terms/history/#titleT-002 |
| **Has Range:** | http://www.w3.org/2000/01/rdf-schema#Literal |

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

# Replication Data for: Regression Discontinuity with Multiple Running Variables Allowing Partial Effects
Version 1.0

Choi, Jin-Young; Lee, Myoung-Jae, 2018, "Replication Data for: Regression Discontinuity with Multiple Running Variables Allowing Partial Effects", doi:10.7910/DVN/PVM8QV, Harvard Dataverse, V1, UNF:6:nmh8KvGK8KXOIEJaGN1L4Q==

≡ Cite Dataset ▾

Learn about Data Citation Standards.

| | |
|---|---|
| **Description** | Data and code to replicate findings in "Regression Discontinuity with Multiple Running Variables Allowing Partial Effects", (2018-01-14) |
| **Subject** | Social Sciences |

Files    Metadata    Terms    Versions

⬆ Export Metadata ▾

Citation Metadata ▲

| | |
|---|---|
| **Dataset Persistent ID** | doi:10.7910/DVN/PVM8QV |
| **Publication Date** | 2018-01-14 |
| **Title** | Replication Data for: Regression Discontinuity with Multiple Running Variables Allowing Partial Effects |
| **Author** | Choi, Jin-Young (Goethe University Frankfurt) Lee, Myoung-Jae (Korea University) |
| **Contact** | ⓘ Use email button above to contact. Choi, Jin-Young (Goethe University Frankfurt) |
| **Description** | Data and code to replicate findings in "Regression Discontinuity with Multiple Running Variables Allowing Partial Effects", (2018-01-14) |
| **Subject** | Social Sciences |
| **Depositor** | Choi, Jin-Young |
| **Deposit Date** | 2018-01-14 |

```xml
<metadata xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:dc="http://purl.org/dc/elements/1.1/"
xmlns:dcterms="http://purl.org/dc/terms/" xmlns="http://dublincore.org/documents/dcmi-terms/">
<dcterms:title>
Replication Data for: Regression Discontinuity with Multiple Running Variables Allowing Partial Effects
</dcterms:title>
<dcterms:identifier>http://dx.doi.org/10.7910/DVN/PVM6QV</dcterms:identifier>
<dcterms:creator>Choi, Jin-Young</dcterms:creator>
<dcterms:creator>Lee, Myoung-Jae</dcterms:creator>
<dcterms:publisher>Harvard Dataverse</dcterms:publisher>
<dcterms:issued>2018-01-14</dcterms:issued>
<dcterms:modified>2018-01-14T17:01:39Z</dcterms:modified>
<dcterms:description>
Data and code to replicate findings in "Regression Discontinuity with Multiple Running Variables Allowing Partial
Effects".
</dcterms:description>
<dcterms:subject>Social Sciences</dcterms:subject>
<dcterms:contributor>Choi, Jin-Young</dcterms:contributor>
<dcterms:dateSubmitted>2018-01-14</dcterms:dateSubmitted>
<dcterms:license>CC0</dcterms:license>
<dcterms:rights>CC0 Waiver</dcterms:rights>
</metadata>
```

# Structured vs. Unstructured
## Metadata



Machine Readable



Human Readable

# Structured Metadata

1. Is encoded in a machine readable format (xml, json)
2. Is compliant with (follows) a standard schema (dublin core, EML, DDI) **OR** accurately defines it's own schema/

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

| Attribute | Value |
|-----------|-------|
| Name | Masha |
| Eye Color | Blue |

# Encoding Attributes

## XML

<eye_color>blue</eye_color>

## JSON

```
{
  "eye_color": "blue",
}
```

**A metadata schema standard will:**

- Define attributes *(e.g. what do you mean by "creator" in ecology?)*

- Suggest controls of values (e.g. dates = MM-DD-YYYY)

- Define requirements for being "well-formed" (e.g. what fields are absolutely necessary for a valid metadata record?)

- Provide an implementation of the standard in an encoding *(e.g. XML)*

- Provide example use cases that are satisfied by the standard.

# Defining the attributes of Open Data

## Dataset Fields

See the *Further Metadata Field Guidance* section to learn more about the use of each element, including the range of valid entries where appropriate. Consult the field mappings to find the equivalent v1.0, DCAT, Schema.org, and CKAN fields.

| Field | Label | Definition | Required |
|-------|-------|-----------|----------|
| @type | Metadata Type | IRI for the JSON-LD data type. This should be `dcat:Dataset` for each Dataset. | No |
| title | Title | Human-readable name of the asset. Should be in plain English and include sufficient detail to facilitate search and discovery. | Always |
| description | Description | Human-readable description (e.g., an abstract) with sufficient detail to enable a user to quickly understand whether the asset is of interest. | Always |
| keyword | Tags | Tags (or keywords) help users discover your dataset; please include terms that would be used by technical and non-technical users. | Always |
| modified | Last Update | Most recent date on which the dataset was changed, updated or modified. | Always |
| publisher | Publisher | The publishing entity and optionally their parent organization(s). | Always |
| contactPoint | Contact Name and Email | Contact person's name and email for the asset. | Always |
| identifier | Unique Identifier | A unique identifier for the dataset or API as maintained within an Agency catalog or database. | Always |
| accessLevel | Public Access Level | The degree to which this dataset **could** be made publicly-available, *regardless of whether it has been made available*. Choices: public (Data asset is or could be made publicly available to all without restrictions), restricted public (Data asset is available under certain use restrictions), or non-public (Data asset is not available to members of the public). | Always |
| bureauCode[USG] | Bureau Code | Federal agencies, combined agency and bureau code from OMB Circular A-11, Appendix C (PDF, CSV) in the format of `015:11`. | Always |
| programCode[USG] | Program Code | Federal agencies, list the primary program related to this data asset, from the Federal Program Inventory. Use the format of `015:001`. | Always |
| license | License | The license or non-license (i.e. Public Domain) status with which the dataset or API has been published. See Open Licenses for more information. | If-Applicable |
| rights | Rights | This may include information regarding access or restrictions based on privacy, security, or other policies. This should also serve as an explanation for the selected "accessLevel" including instructions for how to access a restricted file, if applicable, or explanation for why a "non-public" or "restricted public" data asset is not "public," if applicable. Text, 255 characters. | If-Applicable |

https://project-open-data.cio.gov/v1.1/schema/

# Instructions on how to use…

| Field # | description |
|---|---|
| Cardinality | (1,1) |
| Required | Yes, always |
| Accepted Values | String |
| Usage Notes | This should be human-readable and understandable to an average person. |
| Example | {"description":"This dataset contains a list of vegetables, including nutrition information and seasonality. Includes details on tomatoes, which are really fruit but considered a vegetable in this dataset."} |

https://project-open-data.cio.gov/v1.1/schema/

http://rd-alliance.github.io/metadata-directory/standards/

# Structured Metadata

1. Is encoded in a machine readable format (xml, json)
2. Is compliant with (follows) a standard schema (dublin core, EML, DDI) **OR** accurately defines it's own schema

Think of a custom schema like a playlist.

Individual tracks made by other people, but arranged by you to meet a particular purpose.

In data curation, we call a playlist an **application profile**.

| QDR Field | Dataverse Label | DDI 2.5 | Datacite 3.1 Export |
|---|---|---|---|
| DOI (auto-generated in DV) | Dataset Global ID | 2.1.1.5 IDNo | Identifier with identifierType="DOI" or "Handle" |
| Title | Title | 2.1.1.1 titl | Title |
| | Subtitle | 2.1.1.1 subTitl | Map to Title with titleType="Subtitle" |
| Alternative Title | Alternative Title | 2.1.1.3 altTitl | Map to Title with titleType="AlternativeTitle" |
| QDR ID | Other ID | 2.1.1.5 IDNo | |
| | Identifier | 2.1.1.5 IDNo | alternate identifier |
| | Agency | 2.1.1.5 IDNo (agency) | alternate identifierType |
| Creator | Author | 2.1.2.1 AuthEnty | Creator |
| Name | Name | 2.1.2.1 AuthEnty | creatorName |
| Title/Institutional Affiliation (DDF) | Affiliation | 2.1.2.1 AuthEnty (affiliation) | affiliation |
| Currently only used informally | Identifier | N/A | nameIdentifier |
| Currently only used informally | Identifier Scheme | N/A | nameIdentifierScheme |
| ? | Contact | N/A | Contributor |
| ? | Contact Name | 2.1.4.2 contact | Map to contributorName with contributorType="ContactPers |
| ? | Contact Affiliation | 2.1.4.2 contact (affiliation) | affiliation |
| ? | E-mail | 2.1.4.2 contact (email) | N/A |
| Hardcode QDR | Dataset Publisher | 2.1.4.1 distrbtr | Publisher |
| Publication Date | Publication Date | 2.1.4.5 distDate (for export) | publicationYear |
| Version (automatic) | Version | 2.1.6.1 version | Version |
| Version Date (automatic) | Version Date | 2.1.6.1 version (date) | Map to Date with dateType="Updated" |
| Description | Description | 2.2.2 abstract | Map to Description with descriptionType="Abstract" |
| | Description Date | 2.2.2 abstract (date) | |
| | Subject | 2.2.1.1 keyword | Subject |
| Subject | Keyword | 2.2.1.1 keyword | Subject |

**https://docs.google.com/spreadsheets/d/1kI1Qtq5JneY0cH6yIwORnCHyj7JluGns8VgvcFde4cg/edit#gid=0**

# What makes a good application profile?

*Balances expressivity vs tractability (see last lecture)

*Is applicable to a broad domain

*Uses standards that are complimentary

# 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)

**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>
```

**<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.
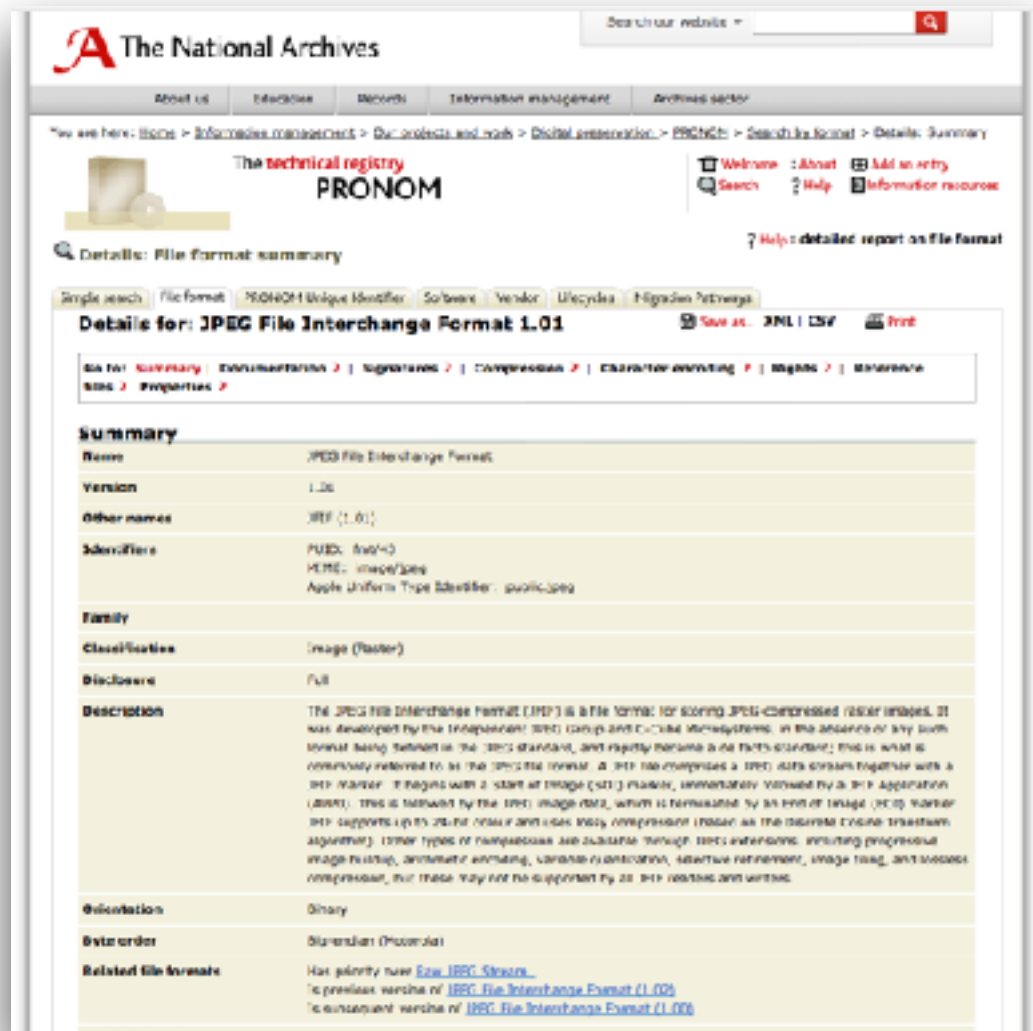
*Example(s):*

<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>

https://www.ddialliance.org/Specification/DDI-Codebook/2.1/DTD/Documentation/titlType.html

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

```xml
<?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 rapidly
became 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>
```

http://www.nationalarchives.gov.uk/PRONOM/Format/proFormatSearch.aspx?status=detailReport&id=668#

# **Unstructured Metadata** or **Documentation**
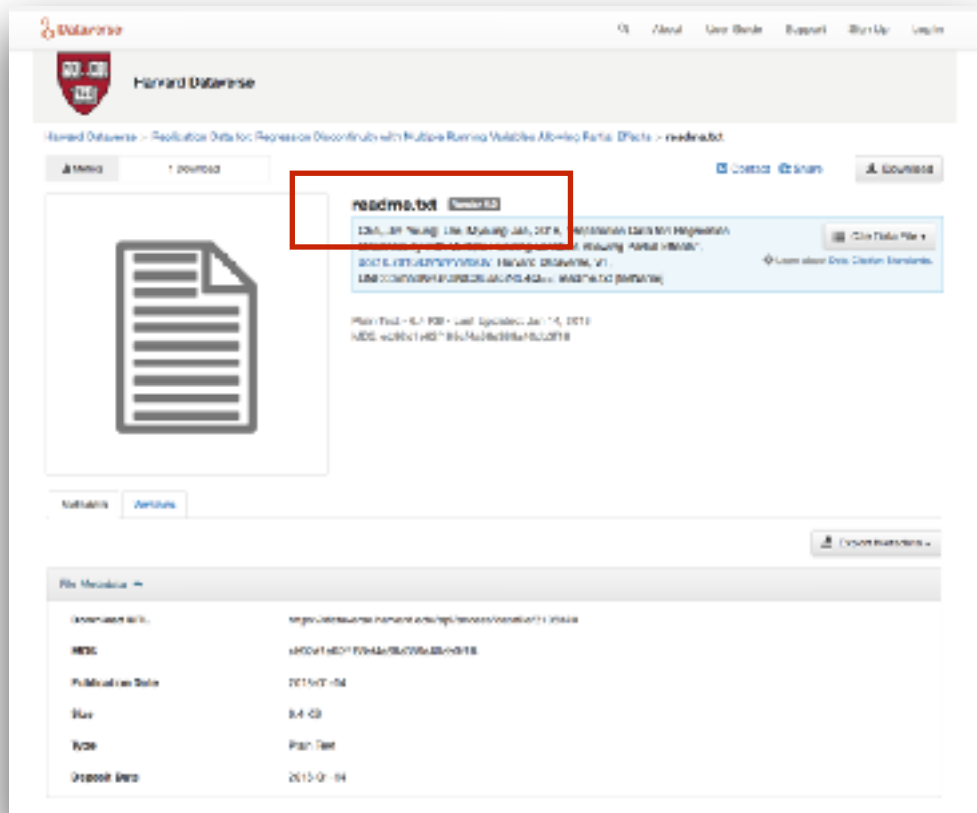(human readable)

**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

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*SD(S)*N^(-1/6)    where S is the running variable in use.

The "mf" value is typically about 0.5-2.5, 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

# Codebook

CODEBOOK FOR ICPSR 9028

UNIFORM CRIME REPORTING PROGRAM DATA [UNITED STATES]

PART 1: OFFENSES KNOWN AND CLEARANCES BY ARREST, 1980

PLEASE NOTE: The "M" between the code and the code label indicates
the code has been designated as a missing value.

| NAME | VARIABLE LABEL | BEG COL | END COL | FMT |
|------|----------------|---------|---------|-----|
| V1 | ID CODE | 1 | 1 | F1 |
| | 1    Offenses known | | | |
| V2 | NUMERIC STATE CODE | 2 | 3 | F2 |

```
1    Alabama
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
```

| V5 | DIVISION | 13 | 13 | F1 |
|----|----------|----|----|----|

```
0    Possessions
1    New England States
2    Middle Atlantic States
3    East North Central States
4    West North Central States
5    South Atlantic States
6    East South Central States
7    West South Central States
8    Mountain States
9    Pacific States
```

| V6 | YEAR | 14 | 17 | F4 |
|----|------|----|----|----|
| V7 | CITY SEQUENCE NUMBER | 18 | 22 | F5 |
| V8 | CORE CITY INDICATION | 23 | 23 | A1 |

```
N    No, not core city of MSA
Y    Yes, core city of MSA
```

| 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
```

https://www.icpsr.umich.edu/icpsrweb/ICPSR/studies/09028

- 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**