# Introduction to Data Science



Last Update: 29 June 2021

# Chapter 2 Data Collection and Acquisition

Papangkorn Inkeaw, PhD



# Outline

Data Collection and Acquisition

- 1. Data Sources
- 2. Data Representation
  - Data Matrix
  - Types of Data
  - Attributes
- 3. Preparing Data
  - Encoding of Categorical Data
  - Normalization and Standardization
  - Data Quality
  - Data Cleaning
    - Inconsistent Datatypes
    - Missing data
    - Duplicate data

## Data Sources



### Questionnaires

- Paper-based questionnaires
- Electronic-based questionnaires
- Online questionnaires



#### Web Servers

Server software, or hardware dedicated to running said software, that can satisfy World Wide Web client requests.



### Web Services

A service offered by an electronic device to another electronic device, communicating with each other via the World Wide Web

## Data Sources



#### Database

An organized collection of data, generally stored and accessed electronically from a computer system

### Logs

- Records of events.
- In computer, for example, a file that records either events that occur in an operating system or other software runs, or messages between different users of a communication software.



### **Online Repositories**

- A <u>repository</u> is a central place in which an aggregation of data is kept and maintained in an organized way, usually in computer storage.
- An <u>online repository</u> is a digital library or archive which is accessible via the internet.

## Data Sources

### **Suggested Data Sources**

- UCI Machine Learning Repository
   <u>https://archive.ics.uci.edu/ml/index.php</u>
- Kaggle https://www.kaggle.com/datasets
- Open Government Data of Thailand
   <u>https://data.go.th/</u>



Data Representation

### Example: Cosmic Dataset

	name	id	align	eye	hair	gender	alive	appearances	first_appear	publisher
	X <sub>1</sub>	$X_2$	<i>X</i> <sub>3</sub>	$X_4$	$X_5$	$X_6$	$X_7$	<i>X</i> 8	<i>X</i> 9	<i>X</i> <sub>10</sub>
	Spider-Man (Peter	Secret	Good	Hazel Eyes	Brown Hair	Male	Living	4043	Aug-62	marvel
<b>x</b> <sub>1</sub>	Parker)						Characters			
	Captain America	Public	Good	Blue Eyes	White Hair	Male	Living	3360	Mar-41	marvel
<b>x</b> <sub>2</sub>	(Steven Rogers)						Characters			
				•••				•••	•••	
	Natalia Romanova	Public	Good	Green Eyes	Red Hair	Female	Living	1050	Apr-64	marvel
$\mathbf{x}_n$	(Earth-616)						Characters			



## Data Matrix

Data Representation

Attributes, Variable, Features, Field,...

$$\begin{array}{cccc} X_1 & X_2 & X_d \\ X_{11} & X_{12} & \cdots & X_{1d} \\ \mathbf{X}_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{array} \begin{array}{c} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \mathbf{X}_n \end{array}$$

 $\mathbf{x}_i$  denotes the *i*th row which is a *d*-tuple given as  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})$ 

 $X_j$  denotes the jth column which is a n-tuple given as  $X_j = (x_{1j}, x_{2j}, \dots, x_{nj})$ 



Data Representation

### Example: Cosmic Dataset

	name	id	align	eye	hair	gender	alive	appearances	first_appear	publisher
	X <sub>1</sub>	<i>X</i> <sub>2</sub>	<i>X</i> <sub>3</sub>	$X_4$	$X_5$	<i>X</i> <sub>6</sub>	$X_7$	<i>X</i> 8	<i>X</i> 9	<i>X</i> <sub>10</sub>
	Spider-Man (Peter	Secret	Good	Hazel Eyes	Brown Hair	Male	Living	4043	Aug-62	marvel
<b>x</b> <sub>1</sub>	Parker)						Characters			
	Captain America	Public	Good	Blue Eyes	White Hair	Male	Living	3360	Mar-41	marvel
<b>x</b> <sub>2</sub>	(Steven Rogers)						Characters			
				•••				•••	•••	
<b>N</b> 7	Natalia Romanova	Public	Good	Green Eyes	Red Hair	Female	Living	1050	Apr-64	marvel
<b>x</b> <sub>n</sub>	(Earth-616)						Characters			



We can write an example  $\mathbf{x}_2$  as

 $\mathbf{x}_2 = (Captain America (Steven Rogers), Public, Good, Blue Eyes, White Hair, Male, Living Characters, 3360, Mar - 41, marvel)$ 

## Types of Data Data Representation



### **Quantitative Data**

- This data can be described using **numbers**.
- **Basic mathematical procedures** are possible on the set.



### **Qualitative Data**

- This data <u>cannot be</u> described using numbers and basic mathematics.
- This data is generally described using natural **categories and language**.

### Data Representation



### Numeric Attributes - Quantitative

- One that has a real-valued or integer-valued domain.
- Such as age, height, grade, frequency, etc.

### Discrete

- Take on a finite or countably infinite set
- Such as integer, grade, number of object, etc.

### Continuous

- Take on any real value
- Such as height, weight, size,
- etc.



### **Categorical Attributes**

- One that has a set-valued domain composed of a set of symbols.
- Such as Gender = {M,F}, Education = {High School, BS, MS, PhD}, etc.

Data Representation

#### Nominal

- Attribute values in the domain are unordered.
- Can only equality (=) compare.
- Such as gender, type of hair, etc.

**Categorical Attributes** 

### Ordinal

- Attribute values are ordered.
- Can both equality (=) and inequality (<, >) compare.
- Such as education, feel (unhappy, OK, happy), etc.

Data Representation



Interval-scaled

- Can compute only differences (addition or subtraction)
- For example, temperature measured in °C or °F.
  - If it is 20 °C on one day and 10 °C on previous day
  - We can talk about a temperature drop of 10°C.
  - We cannot say that it is twice as cold as the previous day.

### Ratio-scaled

- Can compute both differences and ratio between values,
- For example, age.
  - If Jone is 20 years old and Jim is 10 years old.
  - We can say that Jone older than Jim with 10 years.
  - We can say that Jone is twice as old as Jim.

Data Representation

### Summary of data types and scale measures

Provides	Nominal	Ordinal	Interval-scaled	Ratio-scaled
The order of values is known		/	/	/
"Count," aka "Frequency of Distribution"	/	/	/	/
Mode	/	/	/	/
Median		/	/	/
Mean			/	/
Can quantify the difference between each values			/	/
Can add or subtract values			/	/
Can multiple and divide values				/
Has "true zero"				/

https://www.mymarketresearchmethods.com/types-of-data-nominal-ordinal-interval-ratio/



Data Representation

### **Cosmic Dataset**

	name	id	align	eye	hair	gender	alive	appearances	first_appear	publisher
	X <sub>1</sub>	$X_2$	<i>X</i> <sub>3</sub>	$X_4$	$X_5$	<i>X</i> <sub>6</sub>	$X_7$	<i>X</i> 8	<i>X</i> 9	<i>X</i> <sub>10</sub>
	Spider-Man (Peter	Secret	Good	Hazel Eyes	Brown Hair	Male	Living	4043	Aug-62	marvel
<b>x</b> <sub>1</sub>	Parker)						Characters			
	Captain America	Public	Good	Blue Eyes	White Hair	Male	Living	3360	Mar-41	marvel
<b>x</b> <sub>2</sub>	(Steven Rogers)						Characters			
								•••	•••	
**	Natalia Romanova	Public	Good	Green Eyes	Red Hair	Female	Living	1050	Apr-64	marvel
$\mathbf{x}_n$	(Earth-616)						Characters			



# Practice: What is the type of each attribute? Nominal, Ordinal, Interval-scaled or Ratio-scaled

# Encoding of Categorical Data

**Preparing Data** 

- Most of Machine learning algorithms can not handle categorical variables.
- $\rightarrow$  We convert them to numerical values.

### **Nominal variable**

### **One Hot Encoding**

- Map each category to a vector that contains 1 and 0
  - 1 presence of the feature
  - 0 absence of the feature



# Encoding of Categorical Data

**Preparing Data** 

## Ordinal

### **Ordinal Encoding**

- The encoding of variables retains the ordinal nature of the variable
- Each category is <u>assigned a value from 1 through the number of possible values</u> by <u>considering</u> the order of values.



## Encoding of Categorical Data Preparing Data

Align

Bad

Neutral

Good

### Practice

How can we encode the following categorical data?



Hair
Black
Bronze
Brown
Gold
Gray

# Normalization and Standardization

**Preparing Data** 

## Normalization

Scale a variable to have a values between 0 and 1

### **Min-Max Normalization:**

$$x_{normalized} = \frac{x - x_{min}}{x_{max} - x_{min}}$$





# Normalization and Standardization

**Preparing Data** 

## Standardization

Transforms data to have a mean of zero and a standard deviation of 1.

### **Z-score Standardization**

$$x_{standardized} = \frac{x - x}{S.D.}$$





## Data Quality Preparing Data



Source: http://itsadeliverything.com/wordpress/images//accuracyvs-precision.jpg



Data should be:

- Accurate and Precise
- **Complete** Does not have "unknown" or "missing" values
- Consistency Two data items in the data set contradict each other
- Valid Conform to defined business rules or constraints
- **Uniform** Using the same units of measure in all systems
- Unique Does not contain duplicates

## Data Cleaning Preparing Data

**Data Cleaning** is the process of detecting and correcting/removing corrupt or inaccurate records from a record set



# Inconsistent Datatypes

Preparing Data >> Data Cleaning

### We expect that:

Values in a particular attribute must be of a particular datatype, e.g., Boolean, numeric (integer or real), date, etc.

	Values in	1 – Living Characters / 0 – Deceased Characters								
	name	id	appearances	first_appear	publisher					
	<i>X</i> <sub>1</sub>	<i>X</i> <sub>2</sub>	<i>X</i> <sub>3</sub>	<i>X</i> <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>	$X_7$	X <sub>8</sub>	X9	X <sub>10</sub>
<b>x</b> <sub>1</sub>	Spider-Man (Peter Parker)	Secret	Good	Hazel Eyes	Brown Hair	Male	1 /	4043	Aug-62	marvel
<b>x</b> <sub>2</sub>	Captain America (Steven Rogers)	Public	Good	Blue Eyes	White Hair	Male	Living Characters	3360	Mar-41	marvel
$\mathbf{x}_n$	Natalia Romanova (Earth-616)	Public	1	Green Eyes	Red Hair	Female	Living Characters	1050	Apr-64	marvel

# Inconsistent Datatypes

Preparing Data >> Data Cleaning

### How to address the Inconsistent datatypes

- Choose an appropriate datatype
- Transform values in another datatype into the selected datatype

	Values in	/0 – Deceased Characters								
	name	alive	appearances	first_appear	publisher					
	<i>X</i> <sub>1</sub>	<i>X</i> <sub>2</sub>	<i>X</i> <sub>3</sub>	$X_4$	X <sub>5</sub>	<i>X</i> <sub>6</sub>	$X_7$	<i>X</i> <sub>8</sub>	<i>X</i> 9	<i>X</i> <sub>10</sub>
NZ	Spider-Man (Peter	Secret	Good	Hazel Eyes	Brown Hair	Male	Living	4043	Aug-62	marvel
<b>x</b> <sub>1</sub>	Parker)						Characters			
	Captain America	Public	Good	Blue Eyes	White Hair	Male	Living	3360	Mar-41	marvel
<b>x</b> <sub>2</sub>	(Steven Rogers)						Characters			
								•••	•••	
	Natalia Romanova	Public	Good	Green Eyes	Red Hair	Female	Living	1050	Apr-64	marvel
$\mathbf{x}_n$	(Earth-616)		Good				Characters			

1 - Living Characters

### We expect that:

All required measures are known.

	IQ v	Job performance	
	Λ <sub>1</sub>	Λ2	
<b>x</b> <sub>1</sub>	78	NA	
<b>x</b> <sub>2</sub>	84	NA	
<b>X</b> <sub>3</sub>	84	NA	
<b>x</b> <sub>4</sub>	85	NA	
<b>x</b> <sub>5</sub>	99	7	Ī
<b>x</b> <sub>6</sub>	105	10	
<b>X</b> <sub>7</sub>	105	11	
<b>x</b> <sub>8</sub>	106	15	
<b>x</b> 9	108	10	
<b>x</b> <sub>10</sub>	112	10	
<b>x</b> <sub>11</sub>	113	12	
<b>x</b> <sub>12</sub>	115	14	1
<b>x</b> <sub>13</sub>	118	16	1
<b>x</b> <sub>14</sub>	134	12	

Job performances of  $x_1, x_2, x_3$ and  $x_4$  are unknow. They are missing value.

### How to deal with the missing value

**Single Imputation**: Generate a single replacement value for each missing data point.

- Arithmetic Mean Imputation
  - replaces missing values with mean of available values
- Regression Imputation
  - replaces missing values with predicted scores from a regression equation
- Hot-deck Imputation
  - A collection of techniques that impute the missing values with scores from "similar" datapoints, such as nearest neighbor hot-deck and last observation carried forward.
- and etc.



- Compute the arithmetic mean of  $X_2$  from
- Replace the missing values of  $X_2$  by the



# Duplicate Data

Preparing Data >> Data Cleaning

### We expect that:

A data should appear on the dataset one time

	name	id	align	eye	hair	gender	alive	appearances	first_appear	publisher
	<i>X</i> <sub>1</sub>	<i>X</i> <sub>2</sub>	<i>X</i> <sub>3</sub>	$X_4$	X5	X <sub>6</sub>	X <sub>7</sub>	X <sub>8</sub>	<i>X</i> 9	<i>X</i> <sub>10</sub>
<b>N</b> 7	Spider-Man (Peter	Secret	Good	Hazel Eyes	Brown Hair	Male	Living	4043	Aug-62	marvel
<b>x</b> <sub>1</sub>	Parker)						Characters			
	Captain America	Public	Good	Blue Eyes	White Hair	Male	Living	3360	Mar-41	marvel
<b>x</b> <sub>2</sub>	(Steven Rogers)						Characters			
	Spider-Man (Peter	Secret	Good	Hazel Eyes	Black Hair	Male	Living	NA	Aug-62	marvel
<b>x</b> <sub>3</sub>	Parker)						Characters			
	Natalia Romanova	Public	Good	Green Eyes	Red Hair	Female	Living	1050	Apr-64	marvel
<b>x</b> <sub>100</sub>	(Earth-616)			_			Characters			

### We have two recodes of Spider-Man. So, the two recodes are <u>duplicate data</u> Moreover, one contradicts each other

# Duplicate Data

Preparing Data >> Data Cleaning

### How to deal with the duplicate data

- 1. Select one recode that is up-to-date and accurate
- 2. Remove the others

	name	id	align	eye	hair	gender	alive	appearances	first_appear	publisher
	X <sub>1</sub>	$X_2$	<i>X</i> <sub>3</sub>	X_4	X5	X <sub>6</sub>	<i>X</i> <sub>7</sub>	<i>X</i> 8	<i>X</i> 9	<i>X</i> <sub>10</sub>
	Spider-Man (Peter	Secret	Good	Hazel Eyes	Brown Hair	Male	Living	4043	Aug-62	marvel
$\mathbf{x}_1$	Parker)	L	L				Characters			
	Captain America	Public	Good	Blue Eyes	White Hair	Male	Living	3360	Mar-41	marvel
<b>x</b> <sub>2</sub>	(Steven Rogers)						Characters			

					•••		•••			
.,	Natalia Romanova	Public	Good	Green Eyes	Red Hair	Female	Living	1050	Apr-64	marvel
<b>X</b> 100	(Earth-616)						Characters			

We have two recodes of Spider-Man. So, the two recodes are duplicate data

# Practice

### Problem

- จงบอกชนิดข้อมูลของแอตทริบิวต์
  - Eye (Categorical / Numerical)
  - Hair (Categorical / Numerical)
  - Number of appearances
     (Categorical / Numerical)
- จงบอกวิธีการจัดการค่าข้อมูลสูญหาย (Missing Value) <u>ที่เหมาะสม</u>ของแอตทริบิวต์ต่อไปนี้
  - Eye
  - Hair
  - Number of appearances

Name	Fve	Hair	Number of
	LyC	IIGH	appearances
Captain America	Blue	NA	3360
Thor	NA	NA	NA
Benjamin Grimm	Blue	NA	2255
Reed Richards	Brown	Brown	2072
Hulk	Brown	NA	2017
Scott Summers	Brown	Brown	1955
Jonathan Storm	Blue	NA	NA
Robert Drake	Brown	NA	1265
*NA – Missing Value			

# References and Further Study

- Mohammed J. Zaki and Wagner Meira Jr (2014). *Data Mining and Analysis: Fundamental Concepts and Algorithm*. New York, USA: Cambridge University Press.
- Baijayanta Roy (2019). All about Categorical Variable Encoding. url: <u>https://towardsdatascience.com/all-about-categorical-variable-</u> <u>encoding-305f3361fd02</u>.
- Craig K. Enders (2010). *Applied Missing Data Analysis*. New York, USA: Guilford Press.
- Comic Characters (2015). Accessed on 05 October 2019. url: <u>https://github.com/fivethirtyeight/data/tree/master/comic-characters</u>