



PIF
صندوق
الاستثمارات العامة



LEVEL 1: Data Literacy

Day 2

COURSE OUTLINE

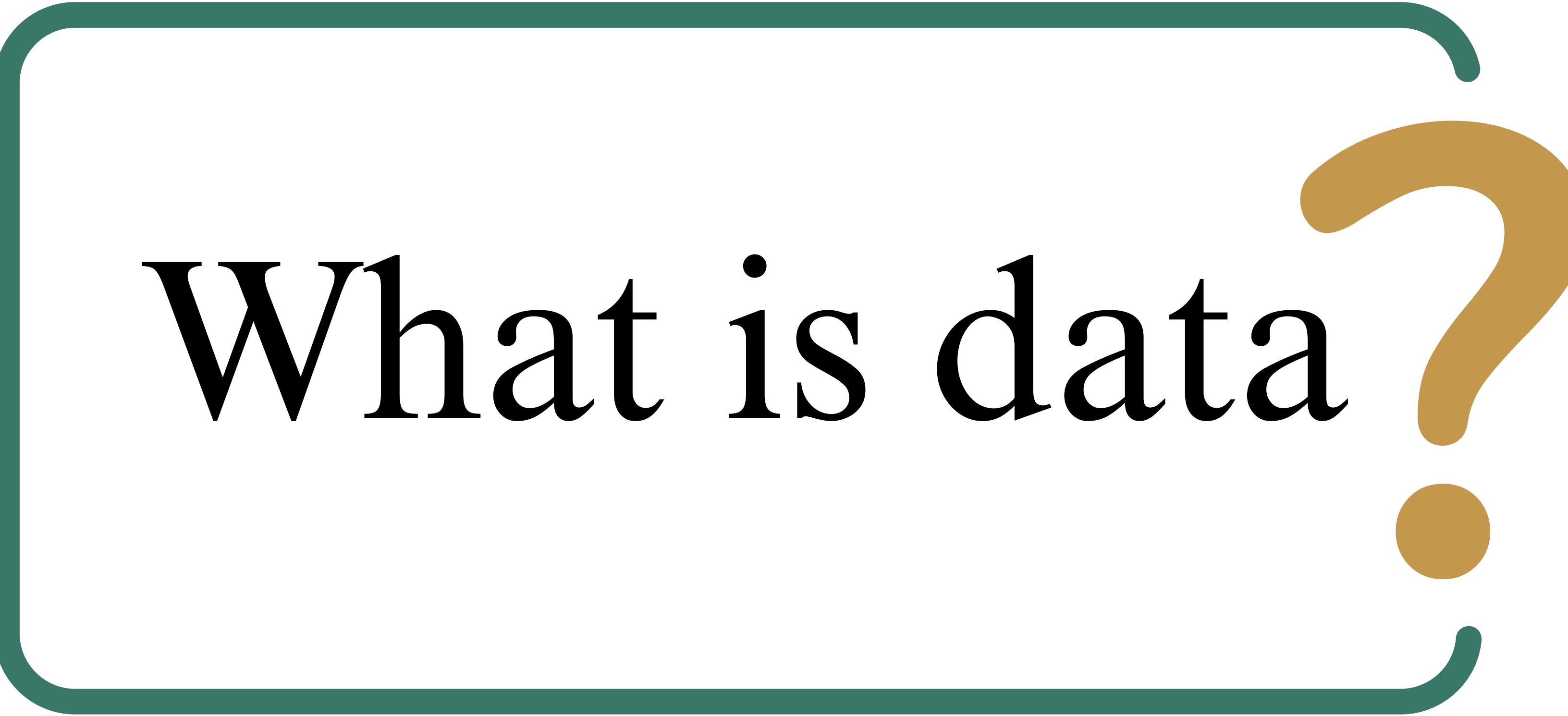
- Definition of Data
- Types of Data
- Data Representation
- Sources of Data
- Data Cleaning

LEARNING OBJECTIVES

- Understand the formal definition of data and its role in analytics and machine learning
- Identify and differentiate between major data types used in ML
- Describe how data can be represented, structured, and encoded for computational use
- Recognize common sources of data across real-world contexts and industries
- Explain the importance of data cleaning and evaluate basic strategies for improving data quality

Data Literacy

Introduction



What is data?

Data Literacy

- Data is a collection of raw facts and observations.



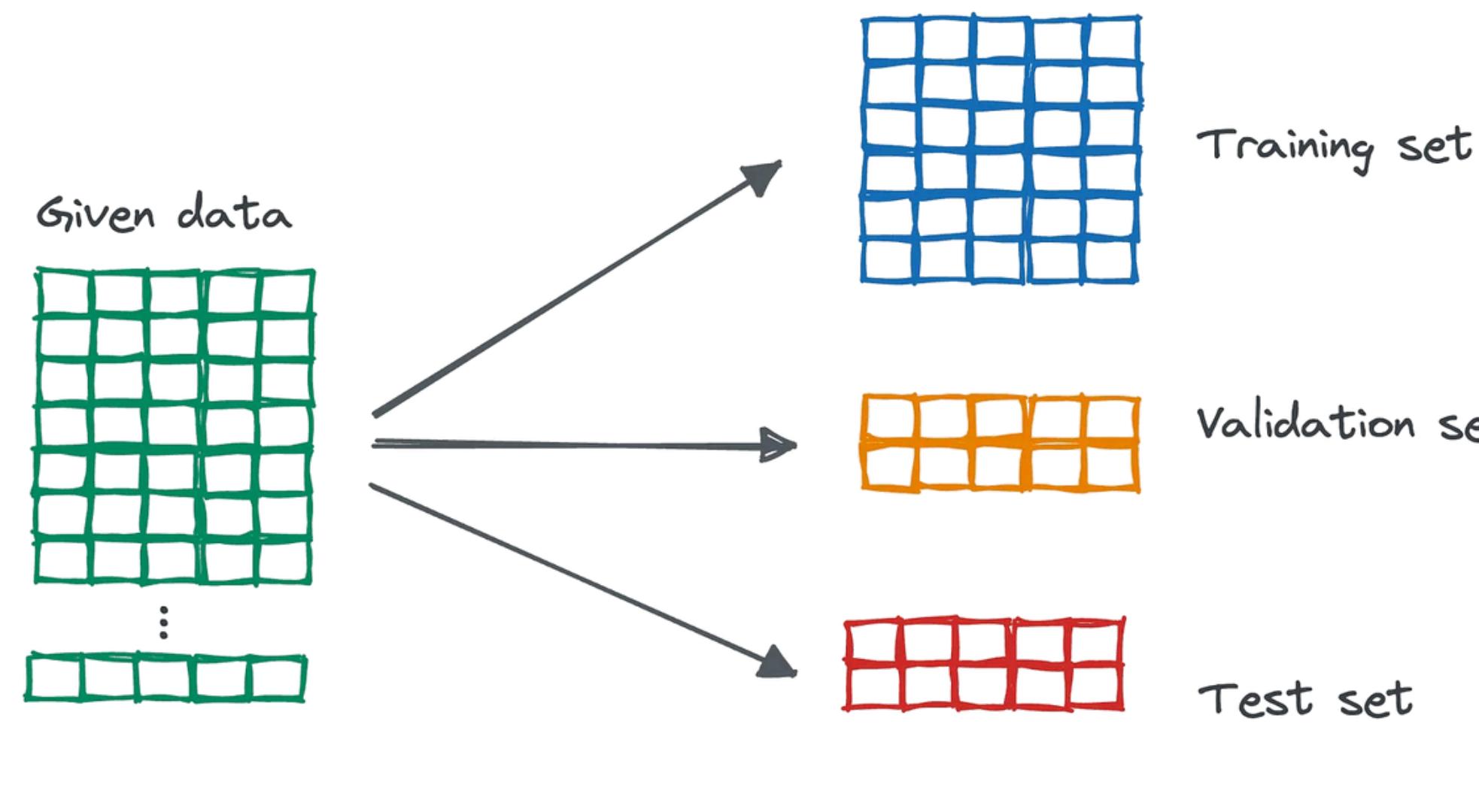
Data Literacy

- Data is the fuel of Machine Learning models.



Data Literacy

- Data used to train, validate, and test a machine-learning model.



Data Literacy

Types of Data

Types of Data

- Tabular Data

columns = attributes for those observations

Rows = observations



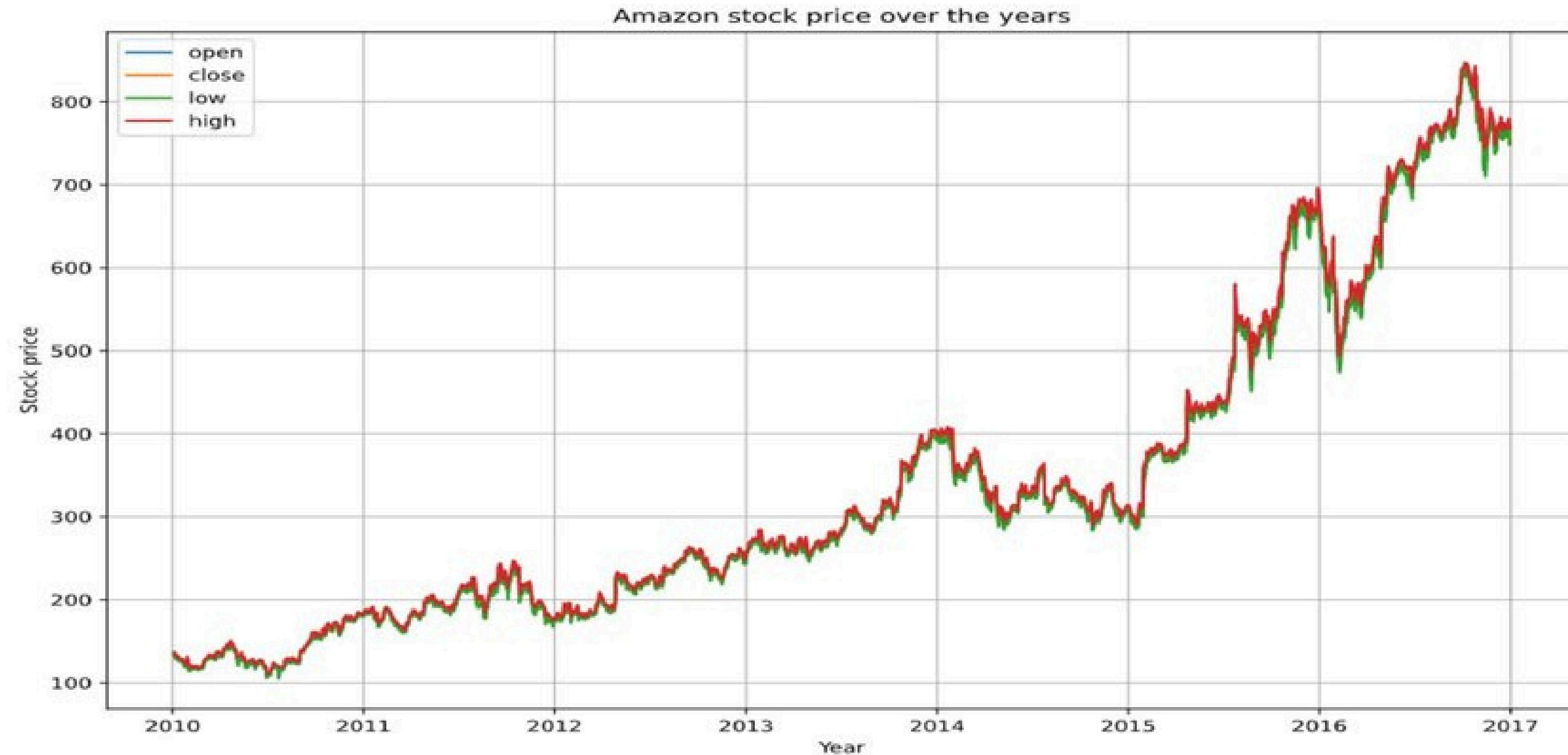
Types of Data

- Tabular Data
 - It is the most common data format in ML.
 - Works well with many traditional ML algorithms.

Size (m ²)	Rooms	Type	Parking	Price (SAR)
120	3	Apartment	1	950,000
180	4	Villa	2	1,350,000
140	3	Apartment	1	980,000
220	5	Villa	3	1,850,000
160	4	Townhouse	2	1,200,000

Types of Data

- Time-Series Data



*

—

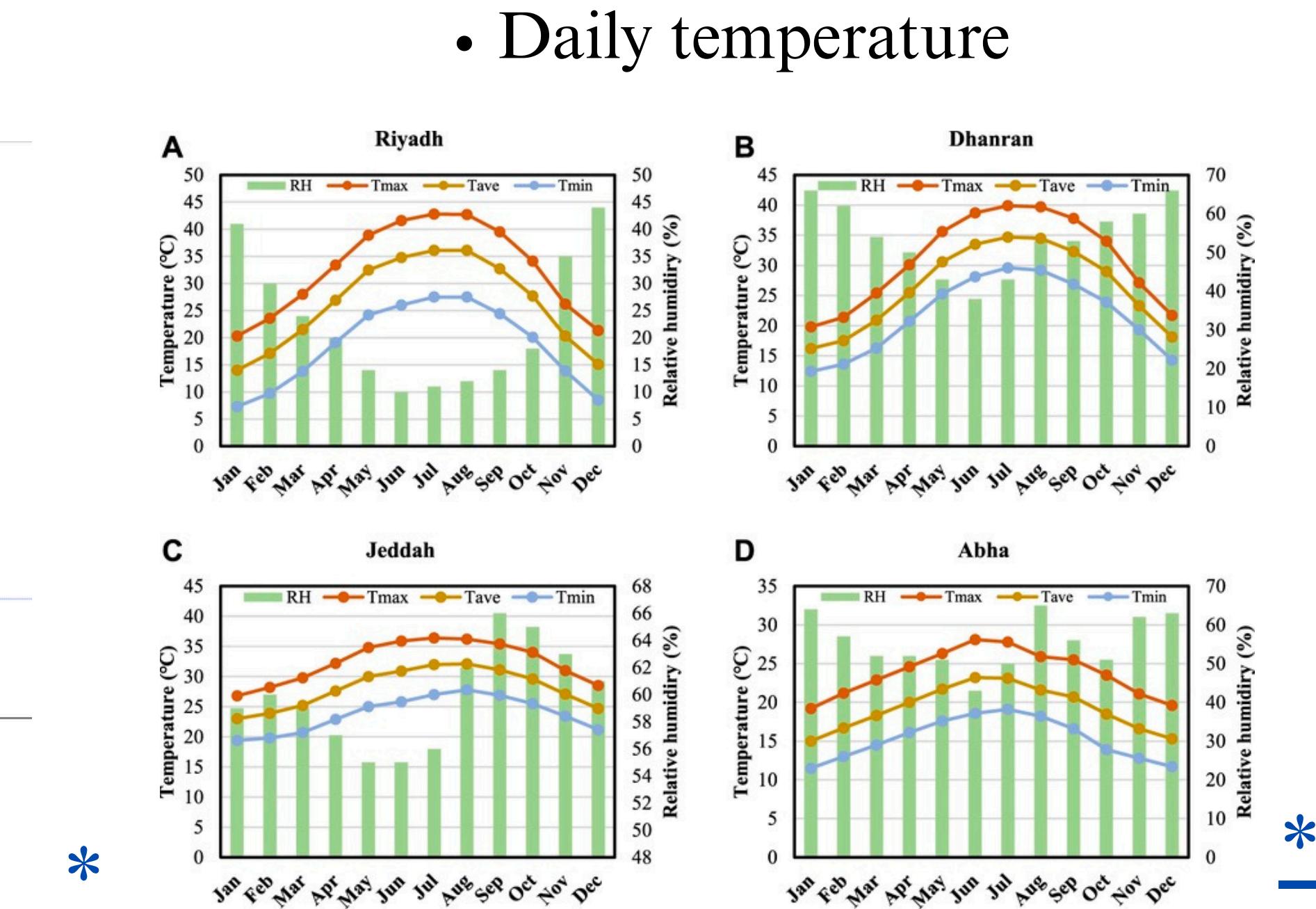
Types of Data

- Time-Series Data
- Data collected over time in sequence.
- Order matters as past affects the future.

Date	Open Price (SAR)	Close Price (SAR)	High	Low	Volume
Jan 1, 2025	128.5	130.1	131.2	127.8	1.2M
Jan 2, 2025	130.15	129.3	131	128.5	1.0M
Jan 3, 2025	129.4	132.6	133.1	129.1	1.5M
Jan 4, 2025	132.7	131.8	134.2	131.1	1.3M
Jan 5, 2025	131.9	135.5	136	131.4	1.6M

Types of Data

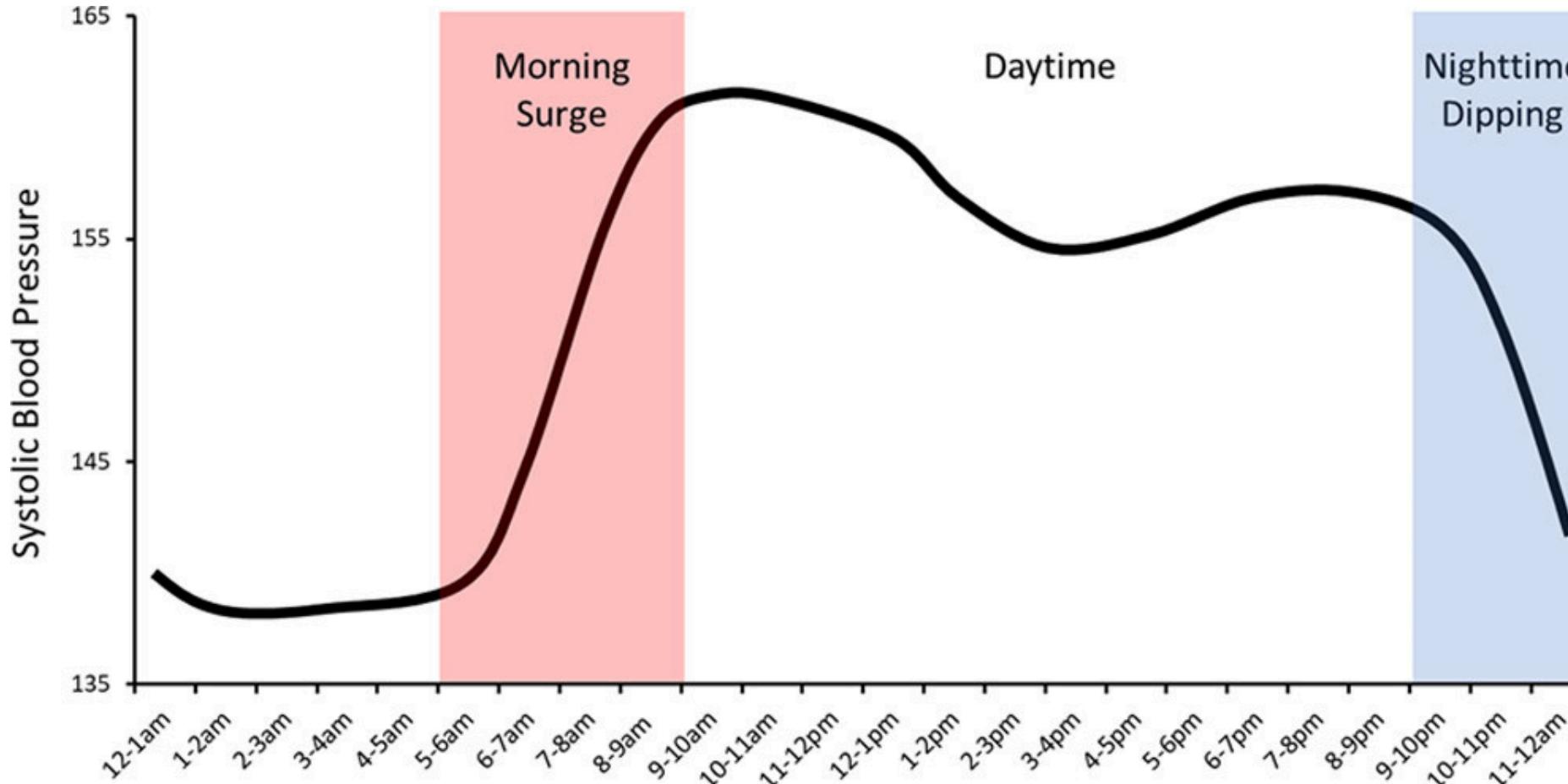
- Time-Series Data
 - Stock-market



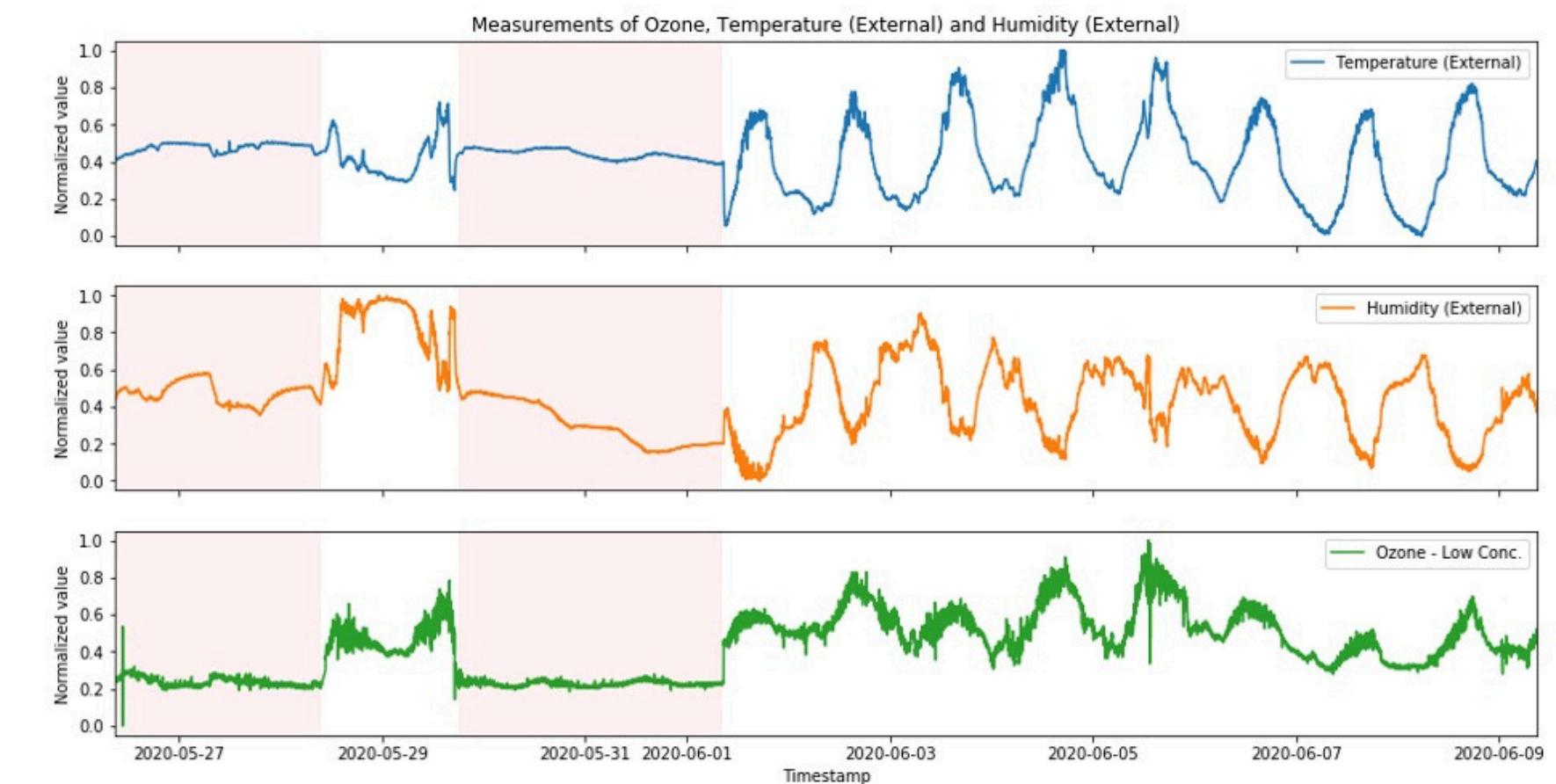
Types of Data

- Time-Series Data

- Blood pressure



- Sensors Data



Types of Data

- Image Data



*

Types of Data

- Image Data

Images are composed of pixels, each encoding specific color and intensity information.

- Colored Images



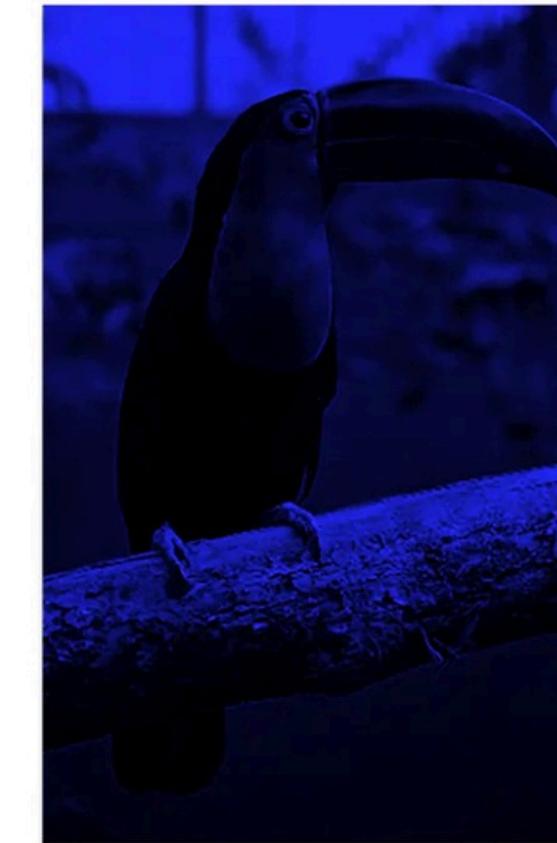
Original Image



Red Channel



Green Channel



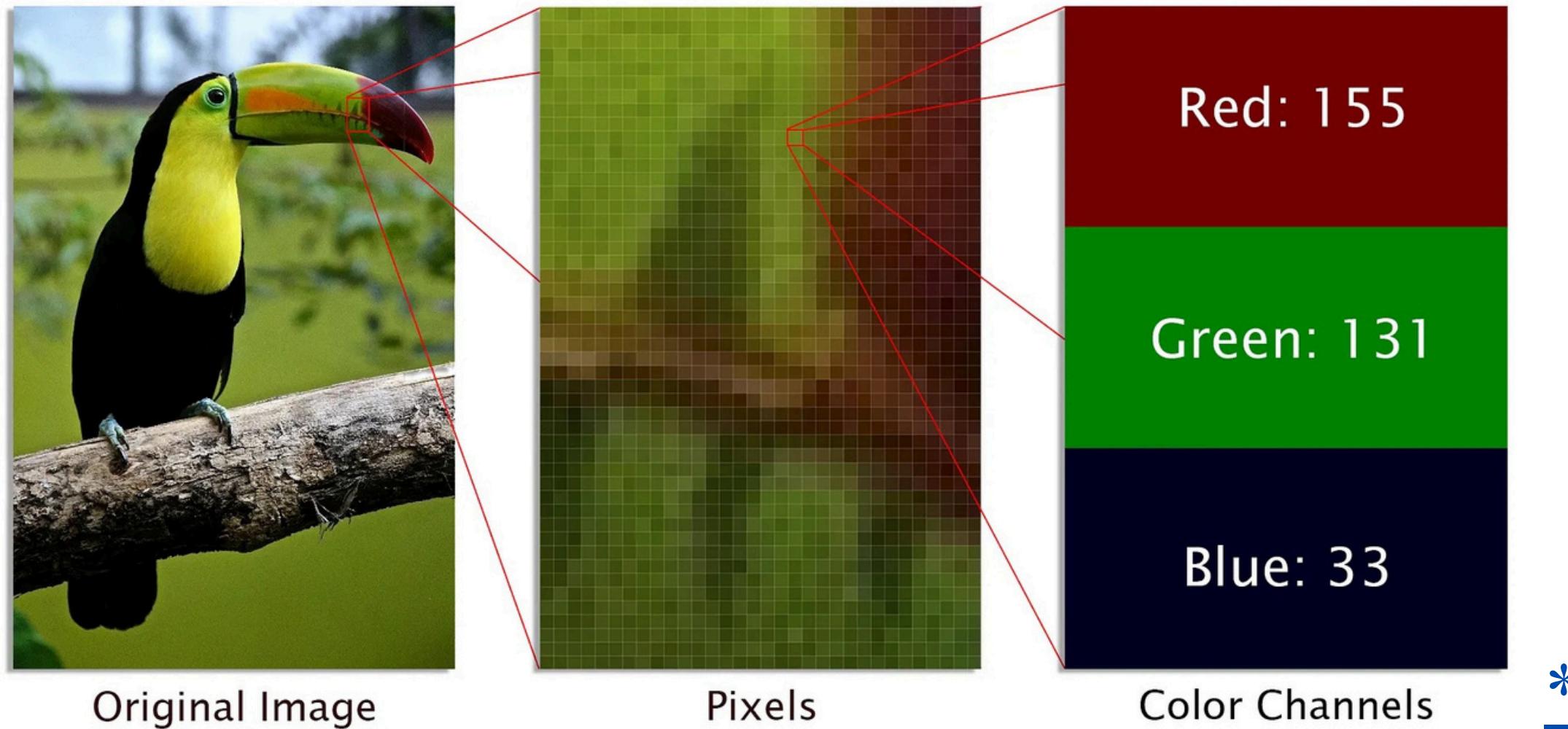
Blue Channel *

Types of Data

- Image Data

Images are composed of pixels, each encoding specific color and intensity information.

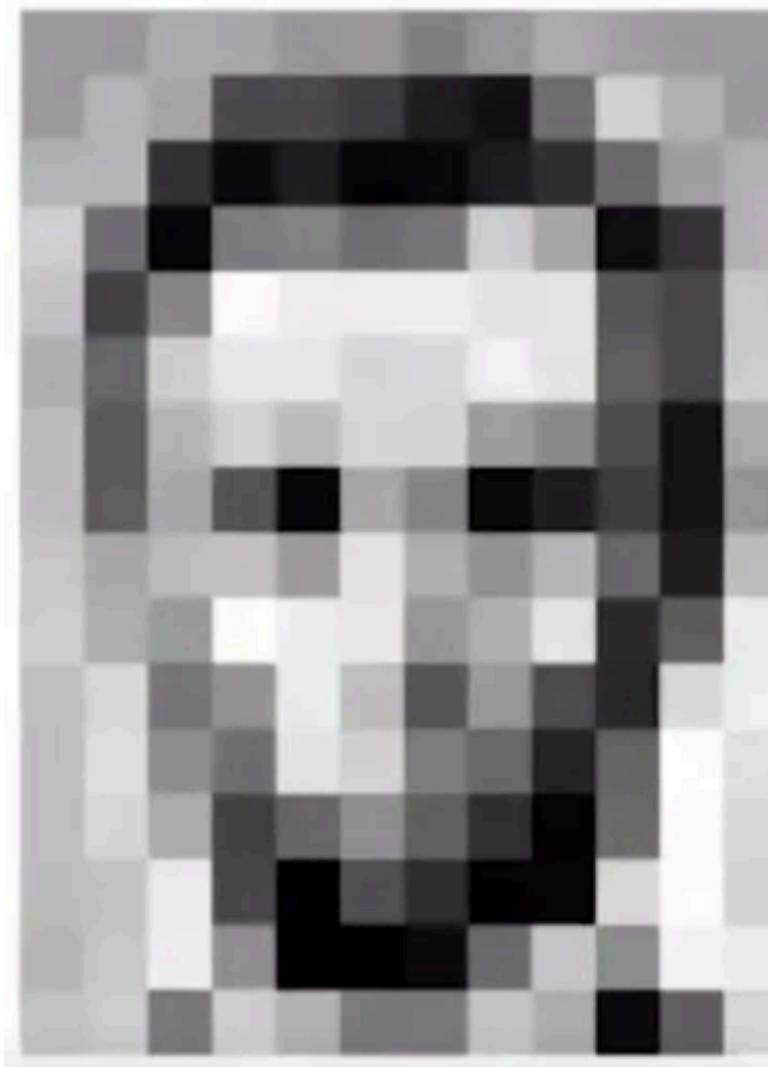
- Colored Images



Types of Data

- Image Data

- Black\white
- Grey images



157	153	174	168	160	162	129	151	172	161	158	156
188	182	163	74	75	62	83	17	112	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	134	107	131	128	204	166	15	54	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	168	219	211	168	138	79	20	169
189	97	168	84	70	168	138	11	31	62	22	148
199	168	191	193	158	227	178	149	182	96	36	190
205	174	155	252	236	231	148	179	238	42	95	234
180	216	119	149	236	187	81	150	79	38	218	241
190	224	147	108	227	210	127	132	94	166	255	224
190	214	173	64	103	143	86	60	7	106	249	215
187	196	238	75	1	81	47	0	6	217	254	211
183	202	237	141	0	0	12	136	200	138	243	236
195	206	123	297	177	135	133	270	178	13	96	218

157	153	174	168	150	152	129	151	172	161	158	156
195	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	168	219	211	168	138	79	20	169
189	97	168	84	70	168	138	11	31	62	22	148
199	168	191	193	158	227	178	149	182	96	36	190
205	174	155	252	236	231	148	179	238	42	95	234
180	216	119	149	236	187	81	150	79	38	218	241
190	224	147	108	227	210	127	132	94	166	255	224
190	214	173	64	103	143	86	60	7	106	249	215
187	196	238	75	1	81	47	0	6	217	254	211
183	202	237	141	0	0	12	136	200	138	243	236
195	206	123	297	177	135	133	270	178	13	96	218

*

Types of Data

Tabular = structured rows & columns

Size (m ²)	Rooms	Type	Parking	Price (SAR)
120	3	Apartmen +	1	950,000
180	4	Villa	2	1,350,000
140	3	Apartmen +	1	980,000
220	5	Villa	3	1,850,000
160	4	Townhous +	2	1,200,000

Image = unstructured pixel grid



260 X 194 X 3

8,11,0, 55,13,25,19
15,241,2,155,13,35,65
14,211,0,255,23,45,11
05,255,1,255,10,17,23
77,167,9,112,56,16,90
45,245,0,145,22,55,48

*

Types of Data

- How to classify these flowers?



Iris Versicolor



Iris Setosa



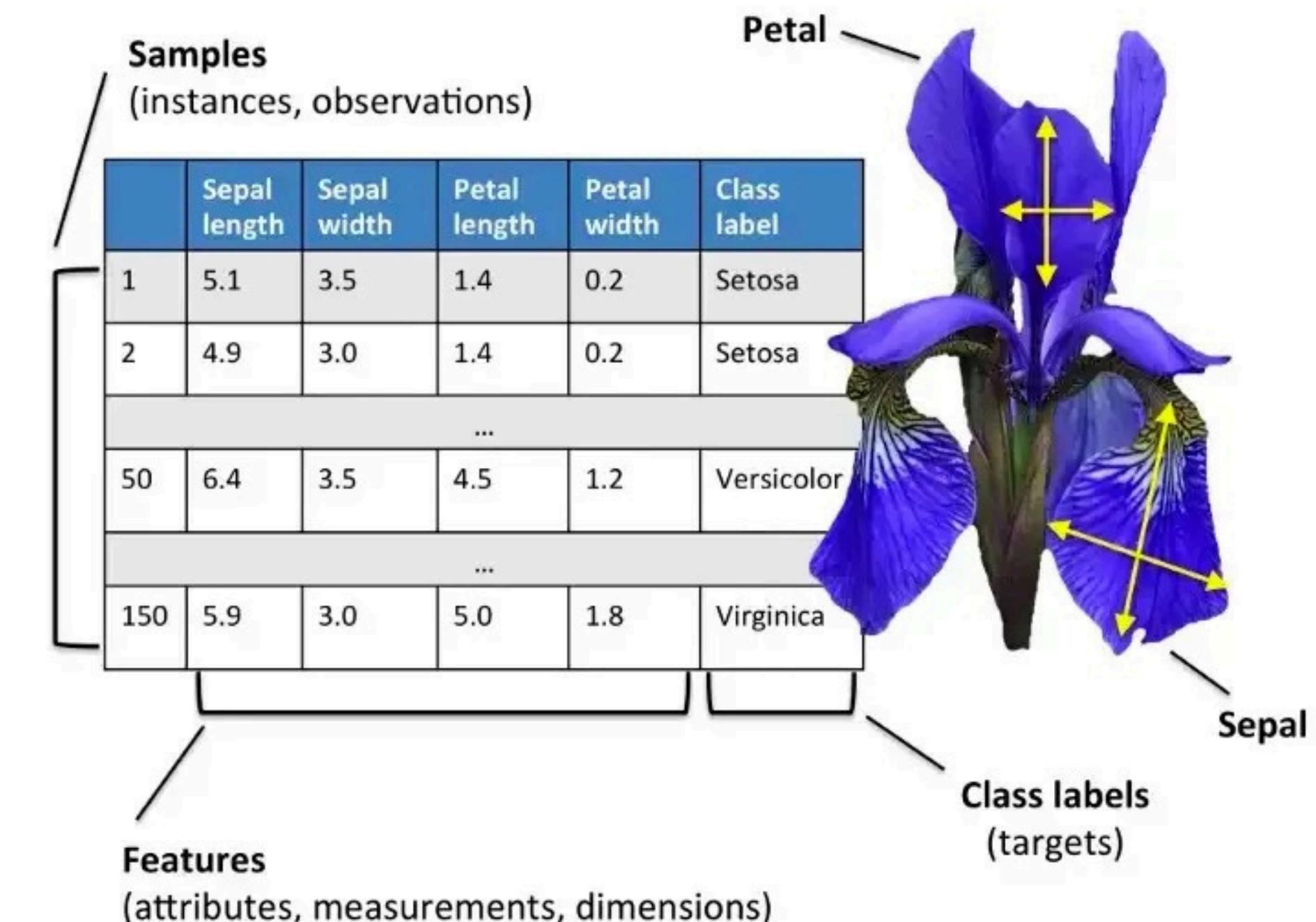
Iris Virginica

*

—

Types of Data

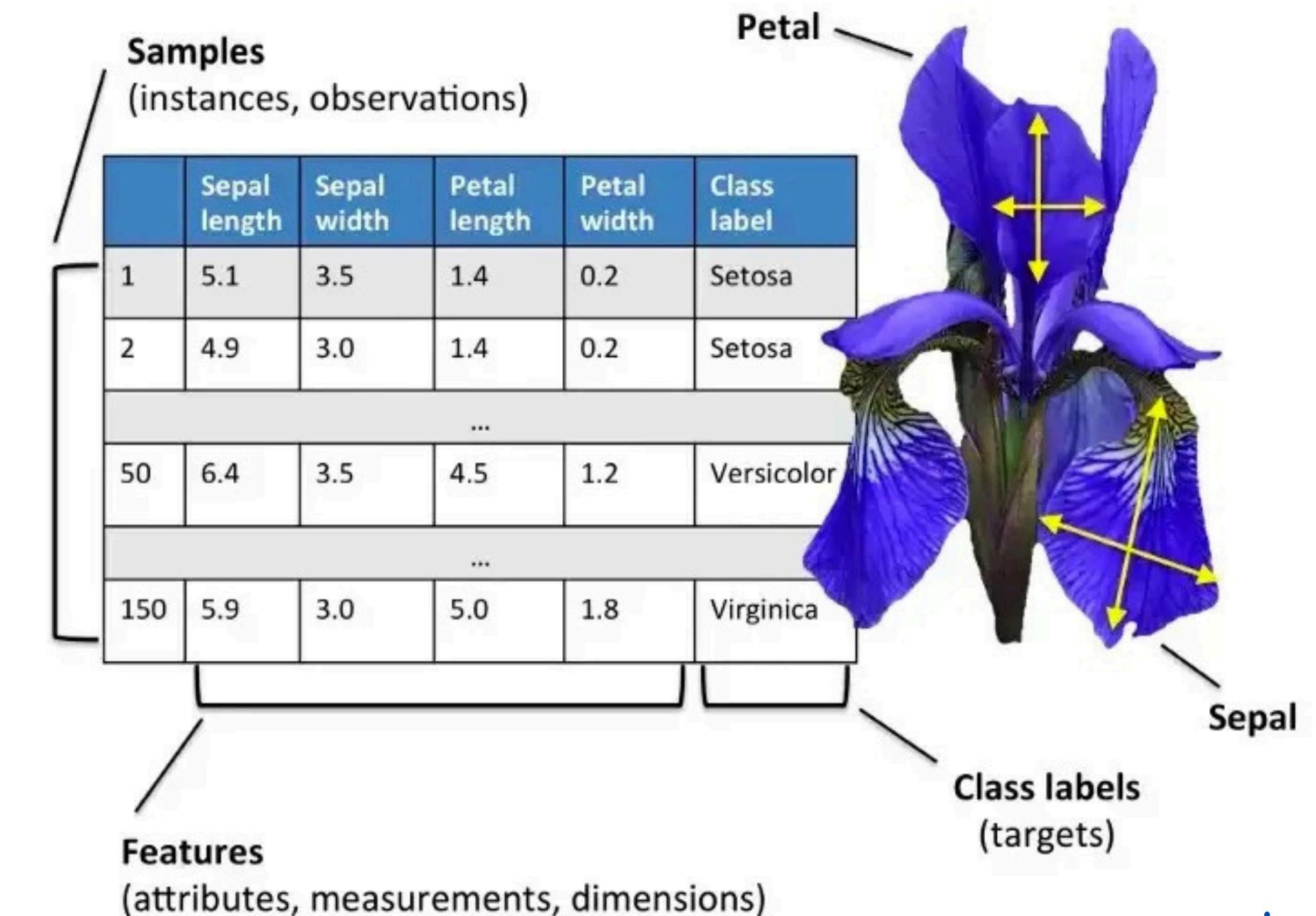
- Traditional tabular data



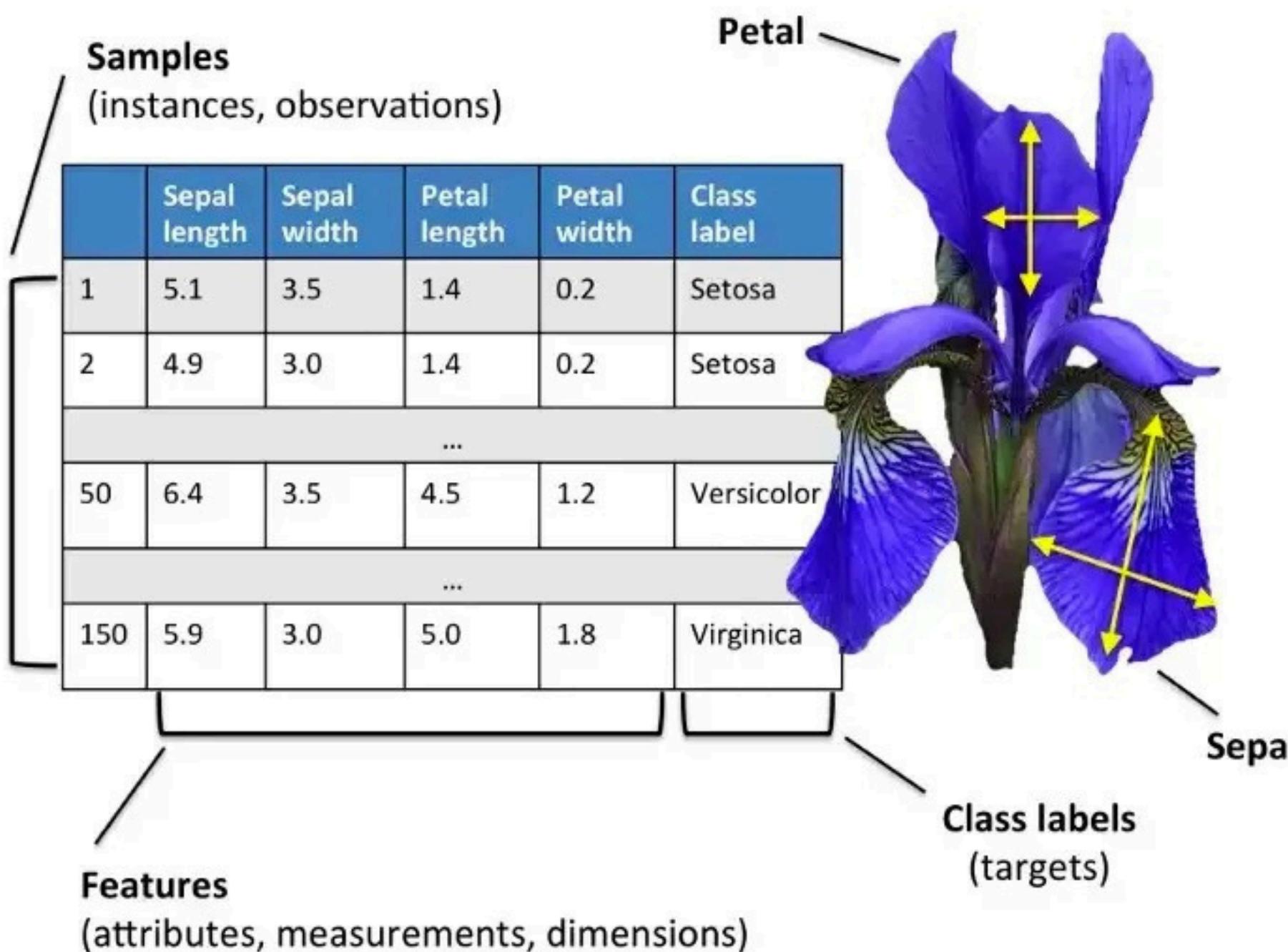
Types of Data

- Traditional tabular data

What's the problem?



Types of Data

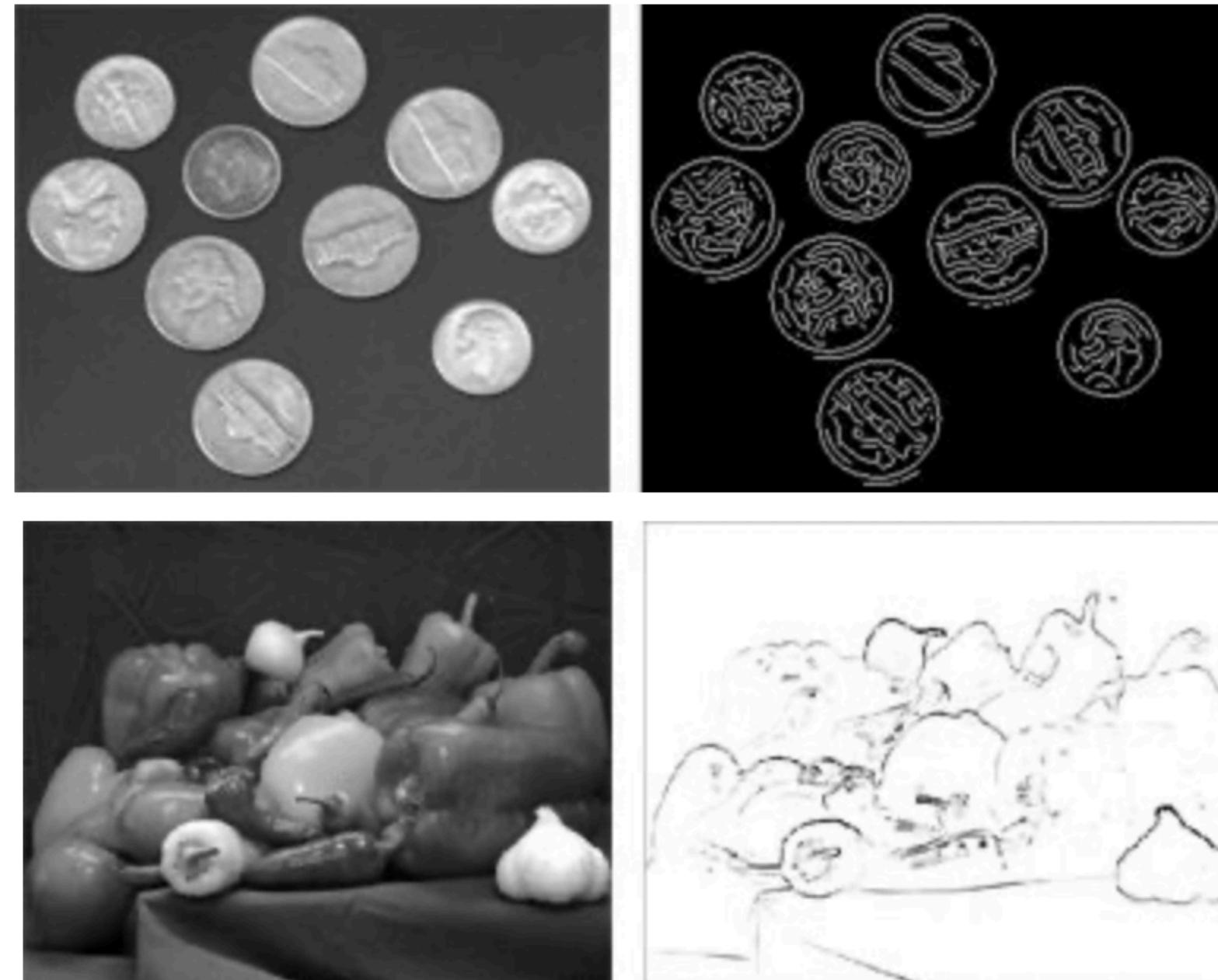


- Human-designed features capture only simple, linear patterns.

Types of Data

- In Machine Learning, images are first transformed into manually engineered features.

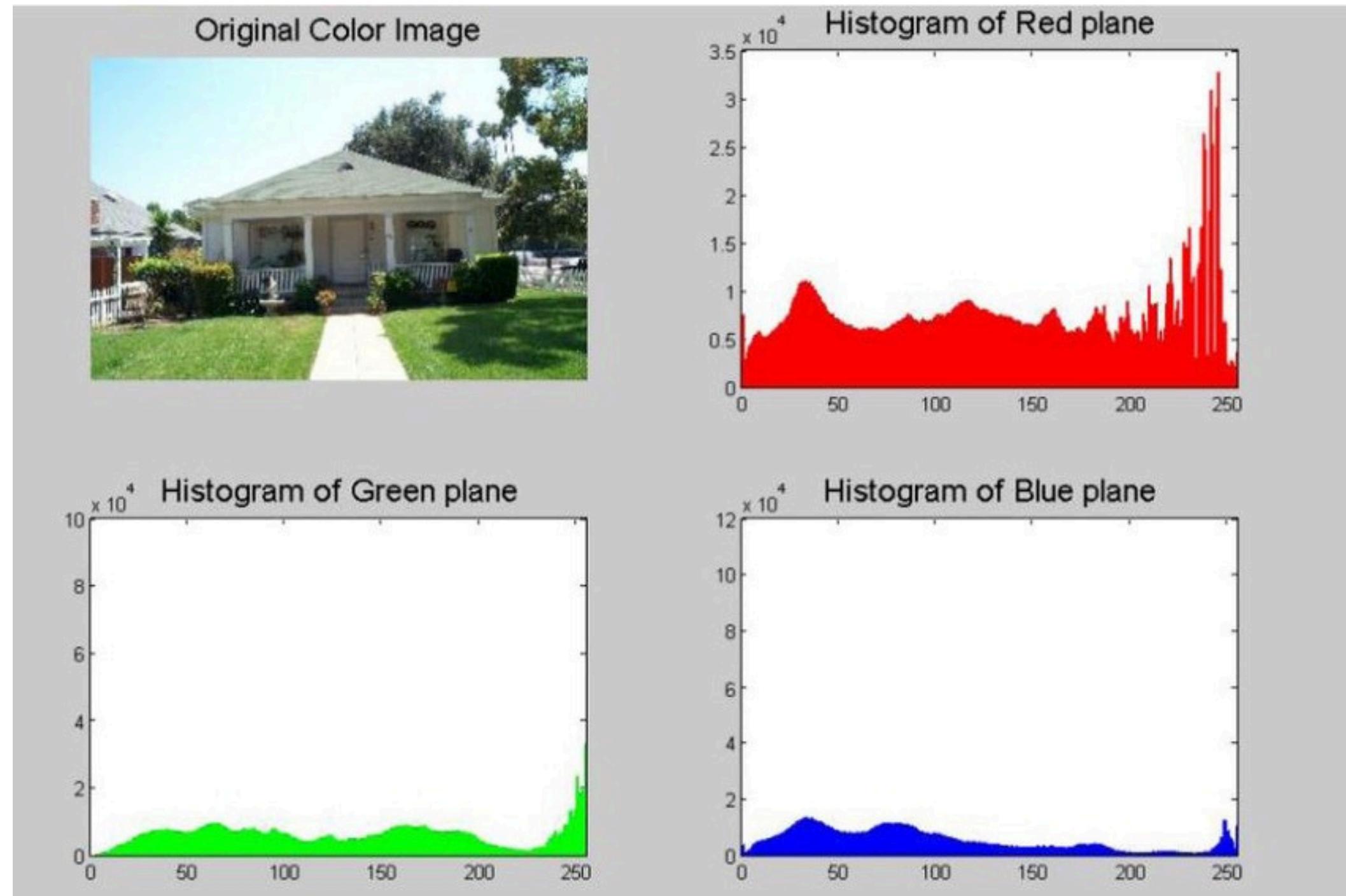
- Edges



*

Types of Data

- Color histograms

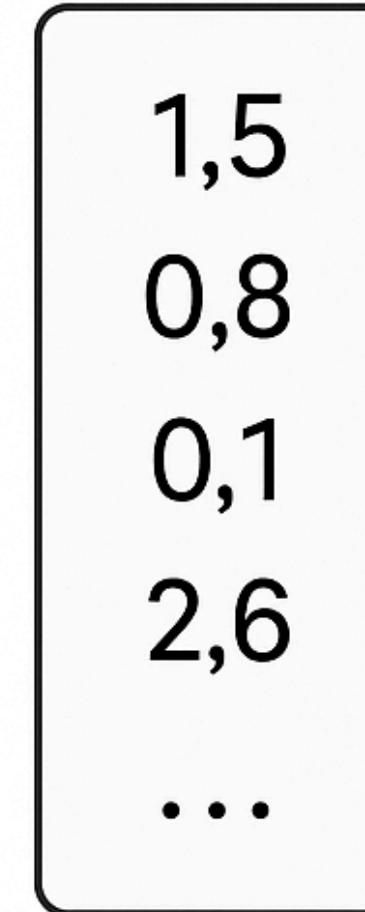


Types of Data



512×512 pixels

Feature
extraction



Feature
vector

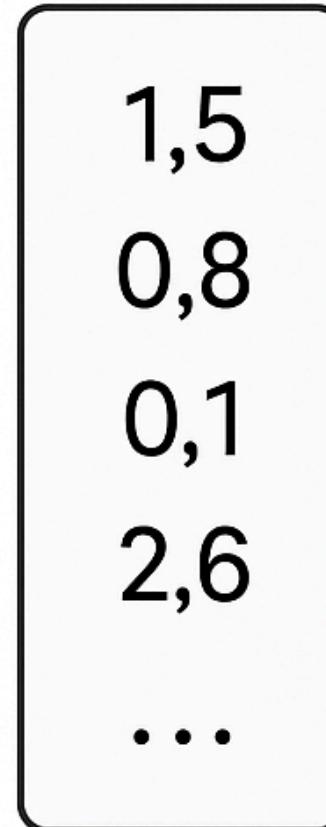
What's the problem here?

Types of Data



512 × 512 pixels

Feature
extraction



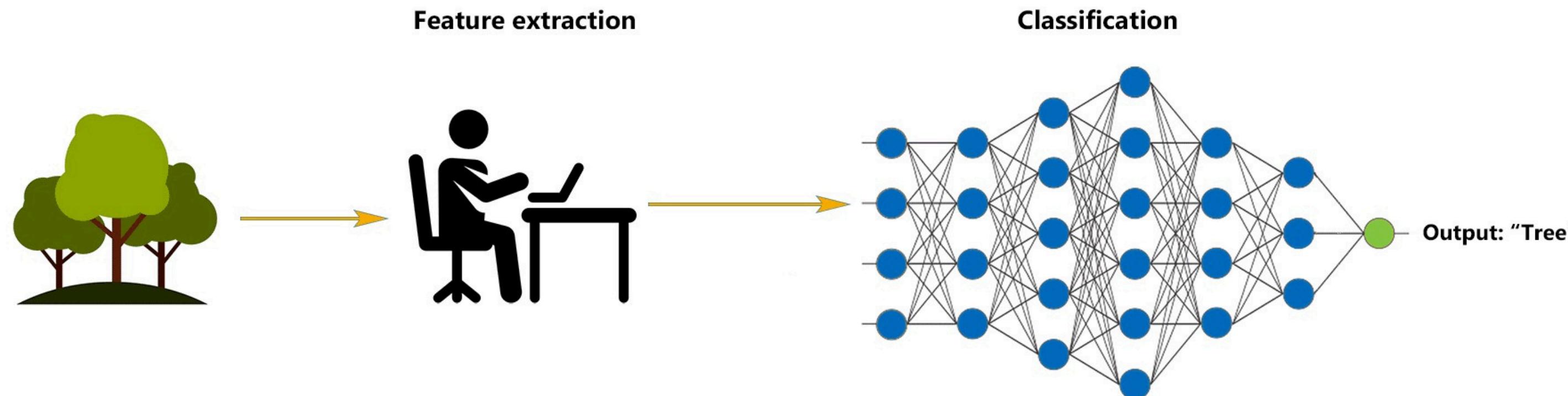
Feature
vector

- Most raw data information is lost when reduced to a few features.

Types of Data

- Traditional ML needs manually extracted features, while deep learning learns features directly from images automatically.

Machine Learning

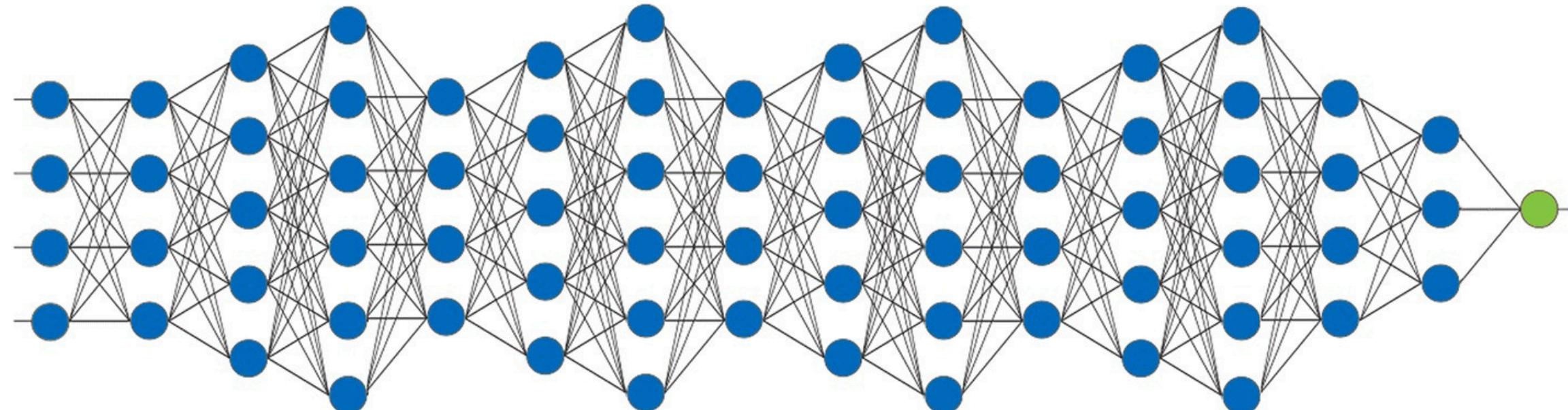


Types of Data

- Traditional ML needs manually extracted features, while deep learning learns features directly from images automatically.

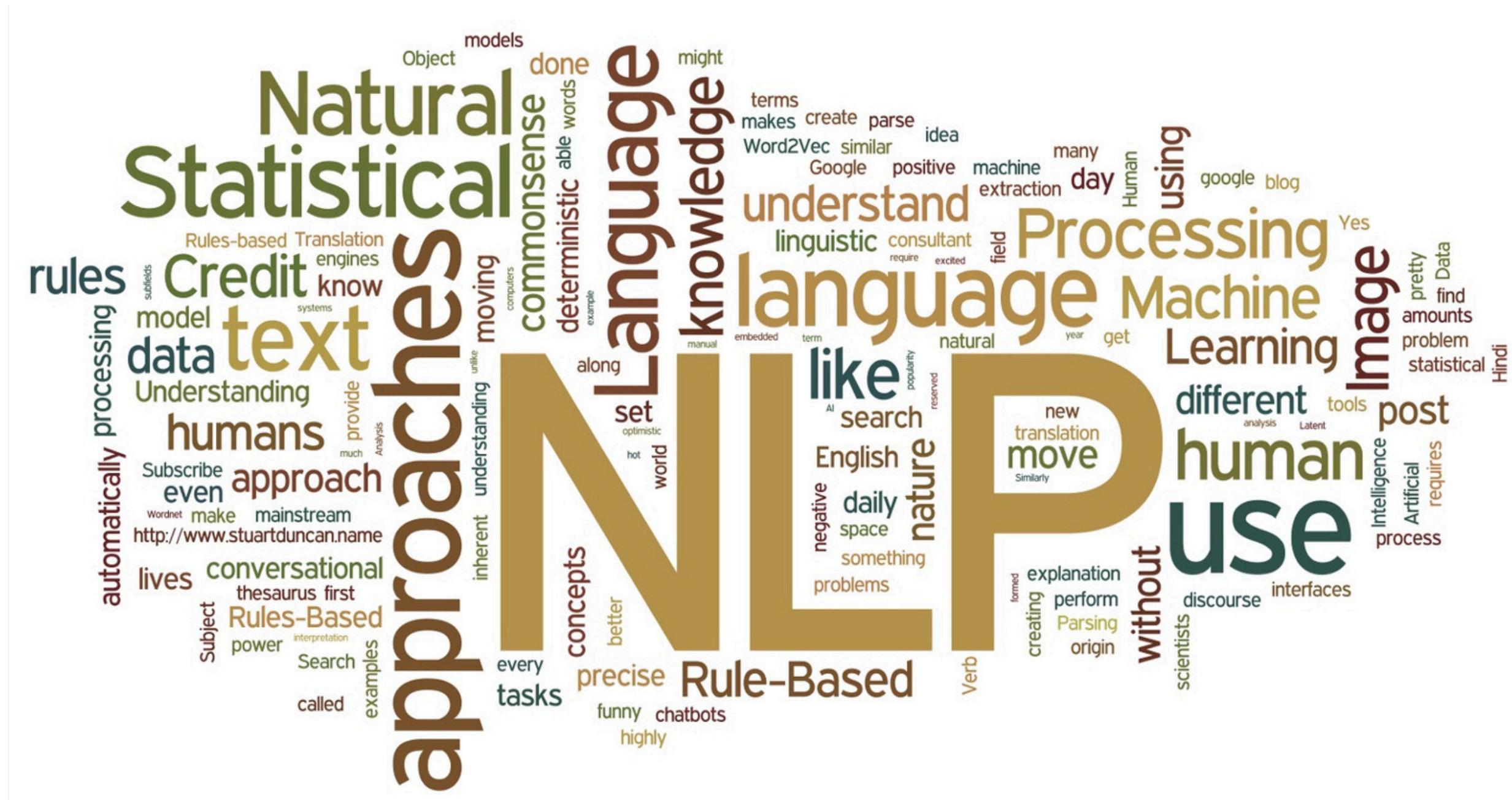
Deep Learning

Feature extraction + Classification



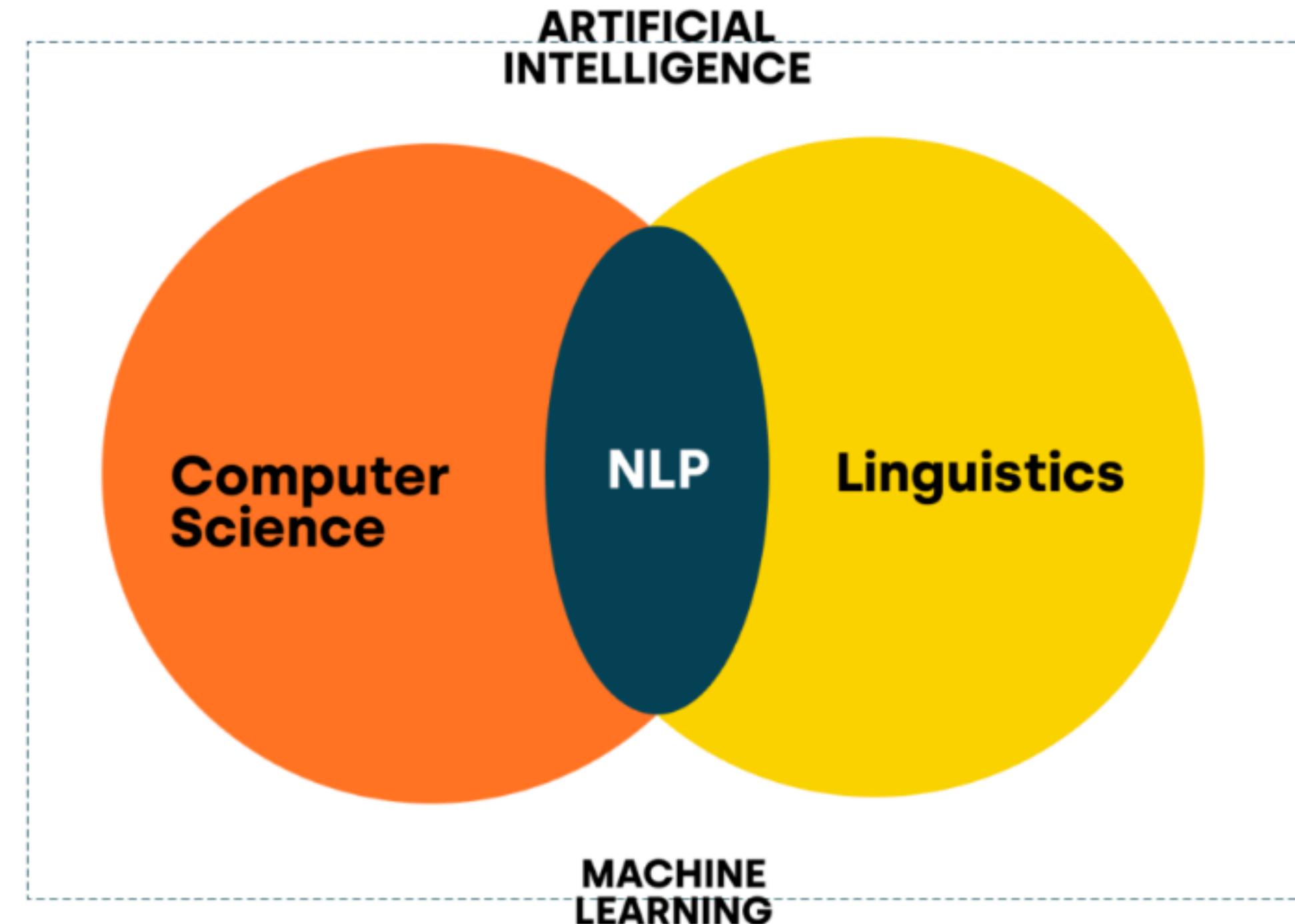
Types of Data

- Textual Data



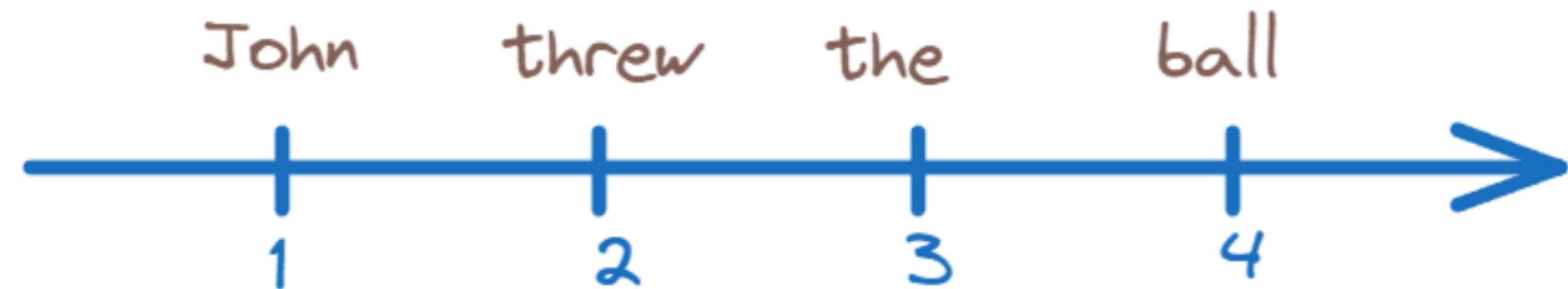
Types of Data

- Textual Data



Types of Data

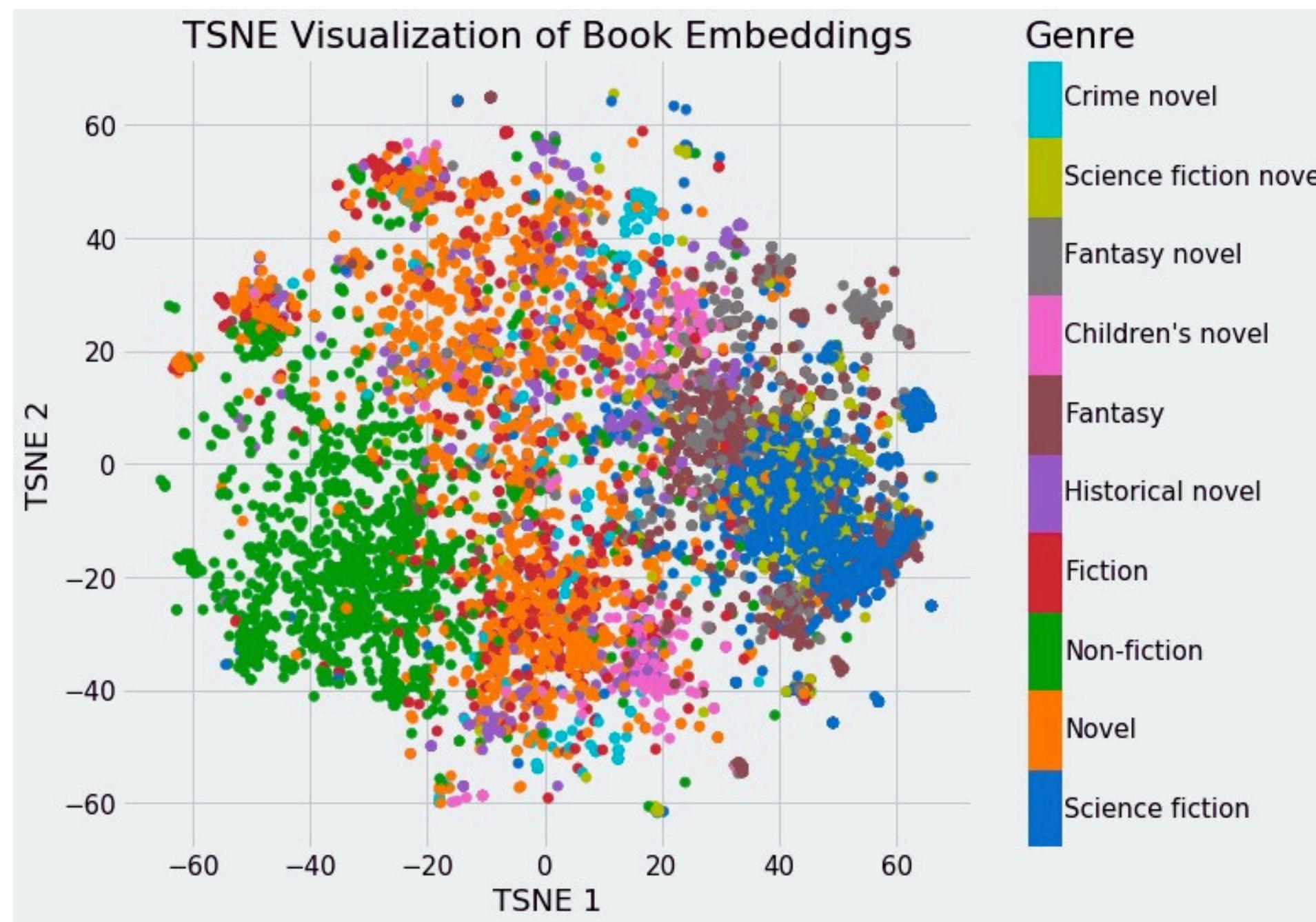
- **Sequential:** words have order and structure.



*

Types of Data

- **High-dimensional:** vocabulary size can be thousands/millions.



Types of Data

- **Ambiguous:** Same word may mean different things in different sentences.



due to its tartness, it is often combined with sweeter juices, such as **apple** or grape



apple is rumored to be working on a smartwatch, which maybe be called an "iwatch."



A cliq app was released for **apple**'s ios devices in august.



word	semantic classes
apple	food, organization

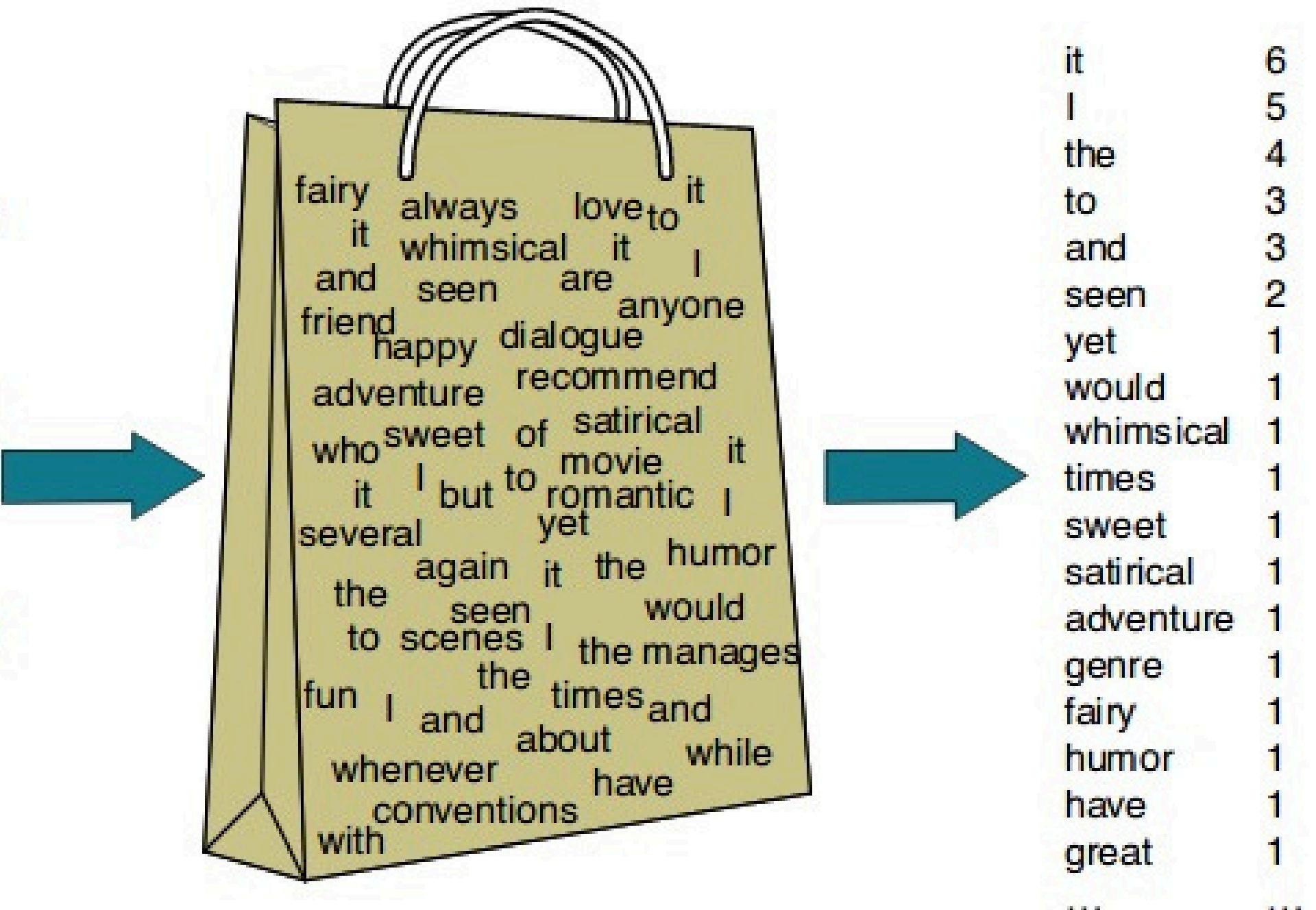


Types of Data

- Text Representation

- Bag of words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



Types of Data

- Text Representation

- Bag of words

Document	NLP	is	fun	very	powerful
Doc1: NLP is fun	?	?	?	?	?
Doc2: NLP is very fun	?	?	?	?	?
Doc3: NLP is powerful	?	?	?	?	?

Types of Data

- Text Representation

- Bag of words

Document	NLP	is	fun	very	powerful
Doc1: NLP is fun	1	1	1	0	0
Doc2: NLP is very fun	1	1	1	1	0
Doc3: NLP is powerful	1	1	0	0	1

Types of Data

- Text Representation
 - TF-IDF
- »»
 - Term Frequency (TF)
 - Measures how often a word appears in a document

$$TF_{ij} = \frac{\text{Number of times the term appears in the document}}{\text{Total number of terms in the document}}$$

Types of Data

- Text Representation
 - TF-IDF
- »»
 - Inverse Document Frequency (IDF)
 - Measures how rare or common a word is across the corpus

$$IDF_i = \log \left(\frac{\text{Total number of documents}}{\text{Number of documents containing term } i} \right)$$

Types of Data

- Text Representation

- TF-IDF

$$TF_{ij} = \frac{\text{Number of times the term appears in the document}}{\text{Total number of terms in the document}}$$

*

$$IDF_i = \log \left(\frac{\text{Total number of documents}}{\text{Number of documents containing term } i} \right)$$

- High IDF \Rightarrow word is rare but valuable
- Low IDF \Rightarrow word is common (e.g., “the”, “a”, “and”)

Types of Data

- Text Representation
 - TF-IDF
- Term Frequency (TF)
 - Measures how often a word appears in a document

$$TF_{ij} = \frac{\text{Number of times term } i \text{ appears in the document}}{\text{Total Number of terms in the document}}$$

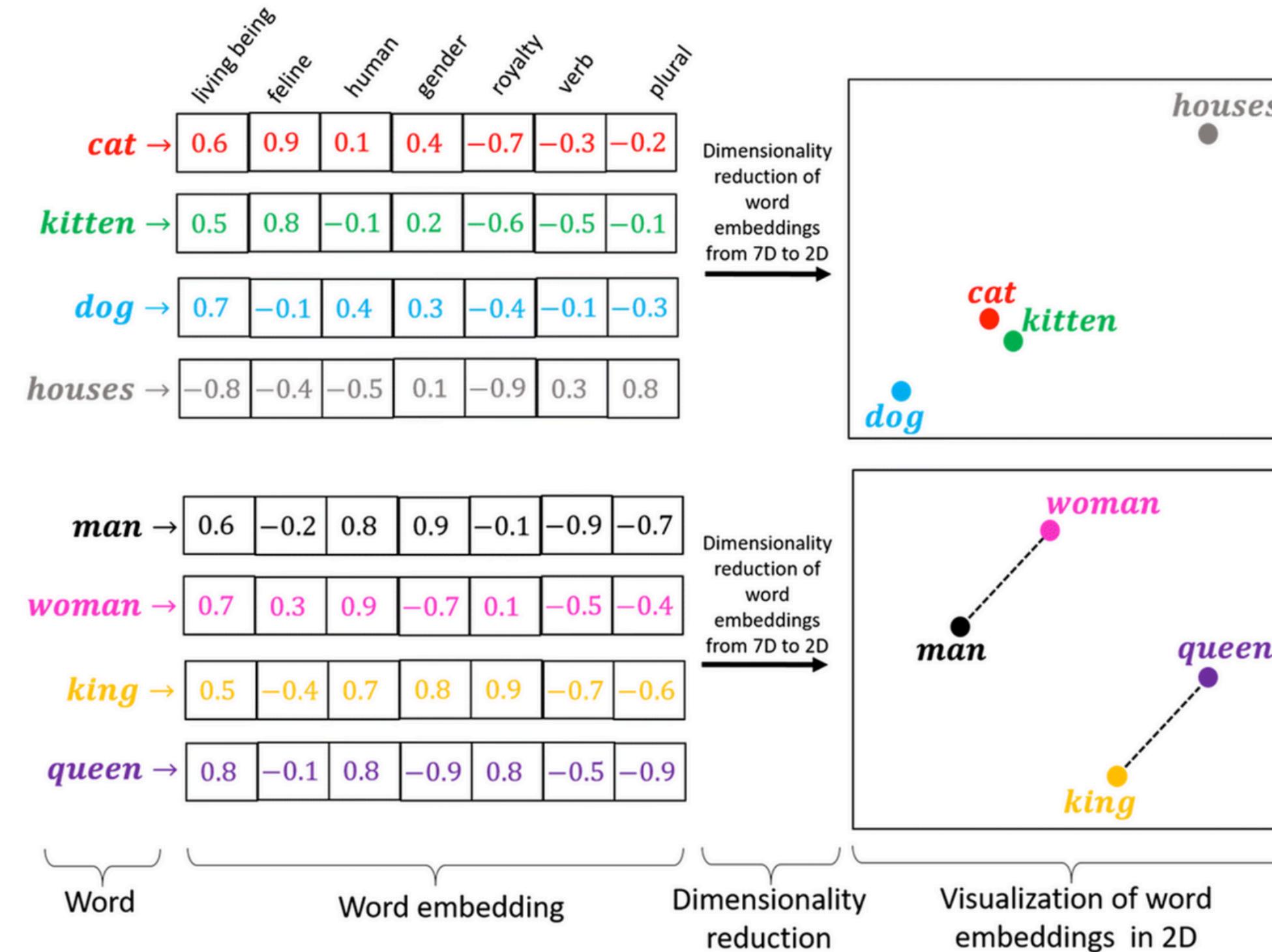
- Inverse Document Frequency (IDF)
- Measures how rare or common a word is across the corpus

$$IDF_i = \log\left(\frac{\text{Total number of documents}}{\text{Number of documents with term } i \text{ in it}}\right)$$

- High IDF \Rightarrow word is rare but valuable
- Low IDF \Rightarrow word is common (e.g., “the”, “a”, “and”)

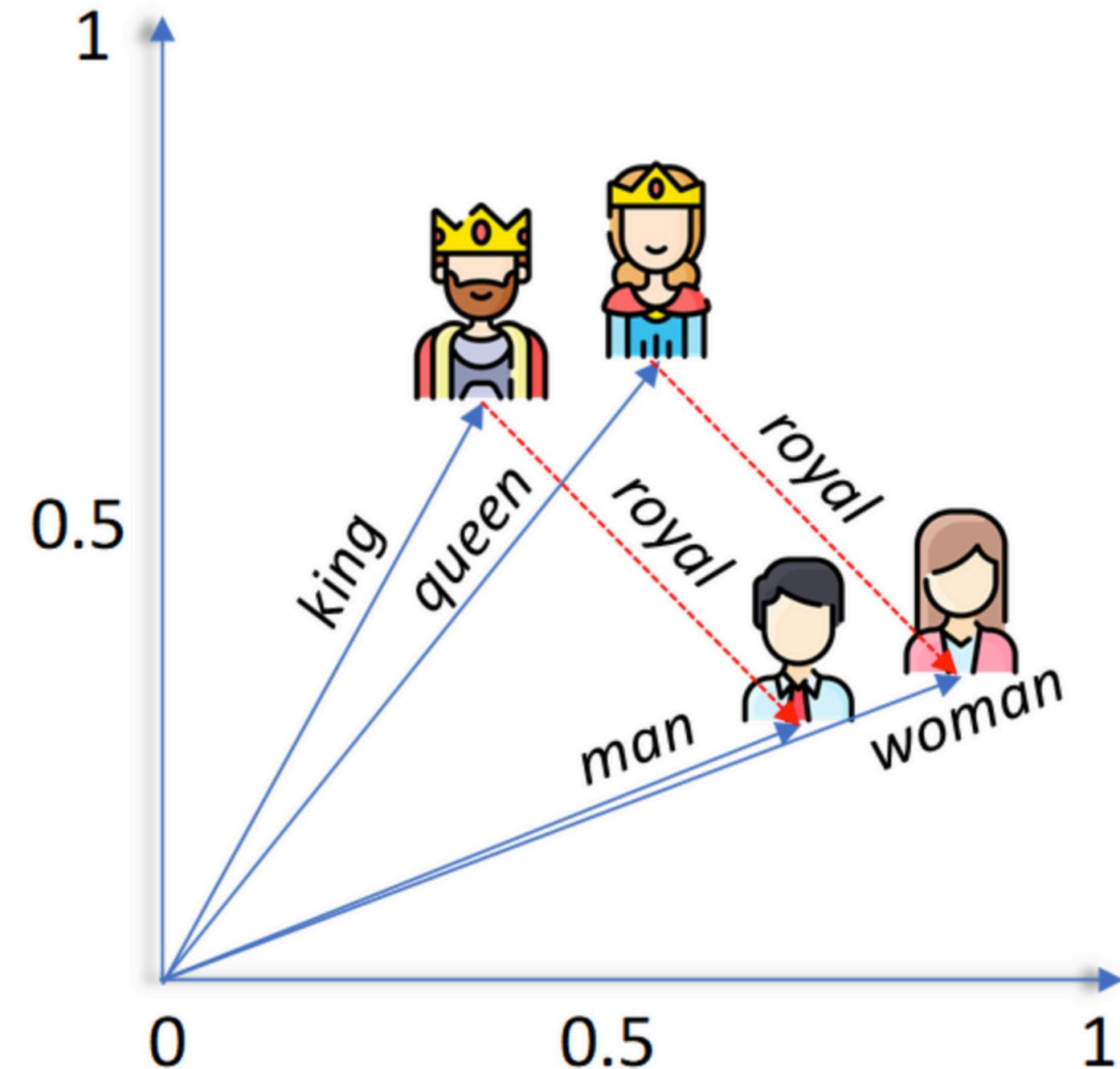
Types of Data

- Word Embedding



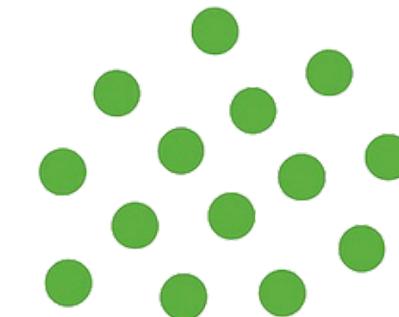
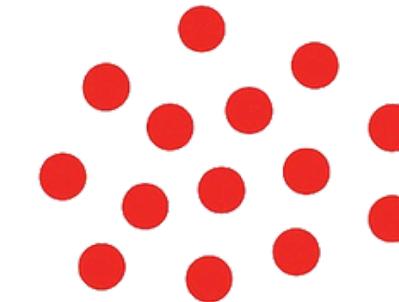
Types of Data

- Word Embedding

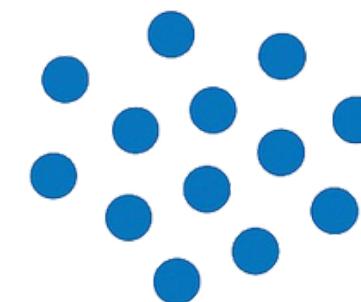


Types of Data

Contextualized embeddings of 'bank'



- Contextualized Embeddings

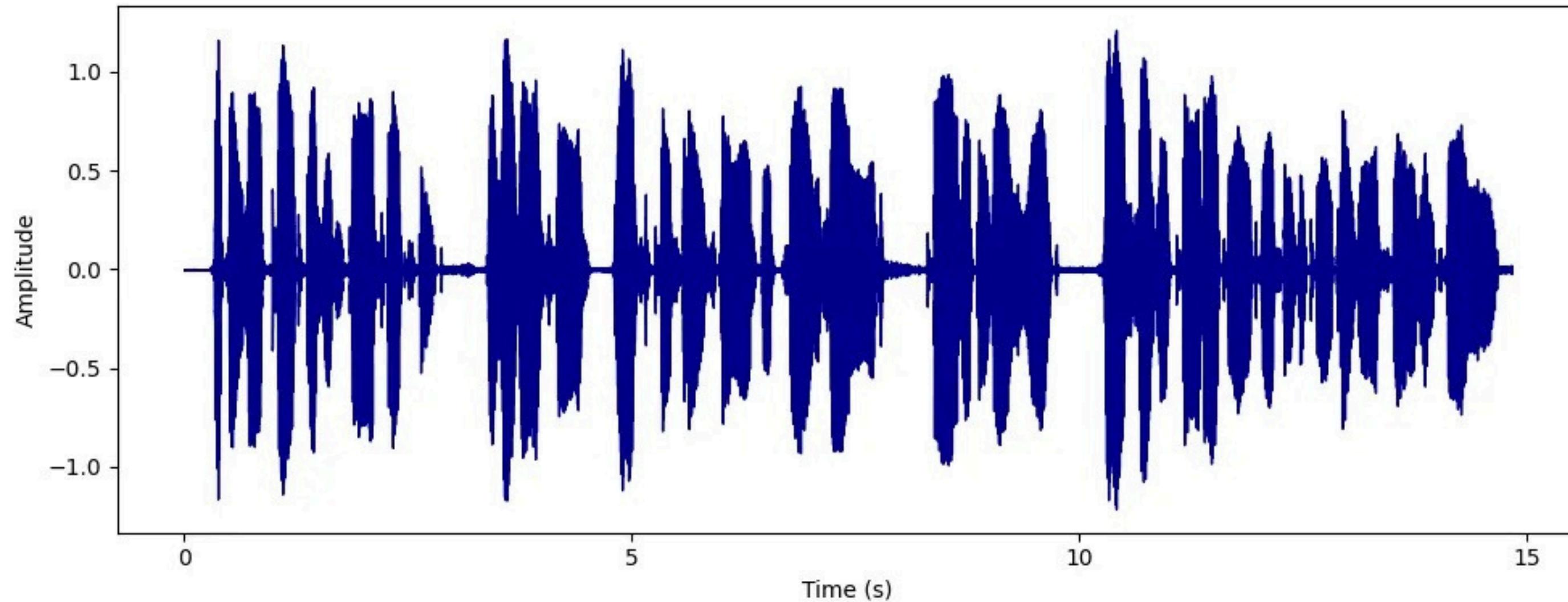


Word
embedding
of 'bank'

bank (financial)
bank (river)
bank (land)

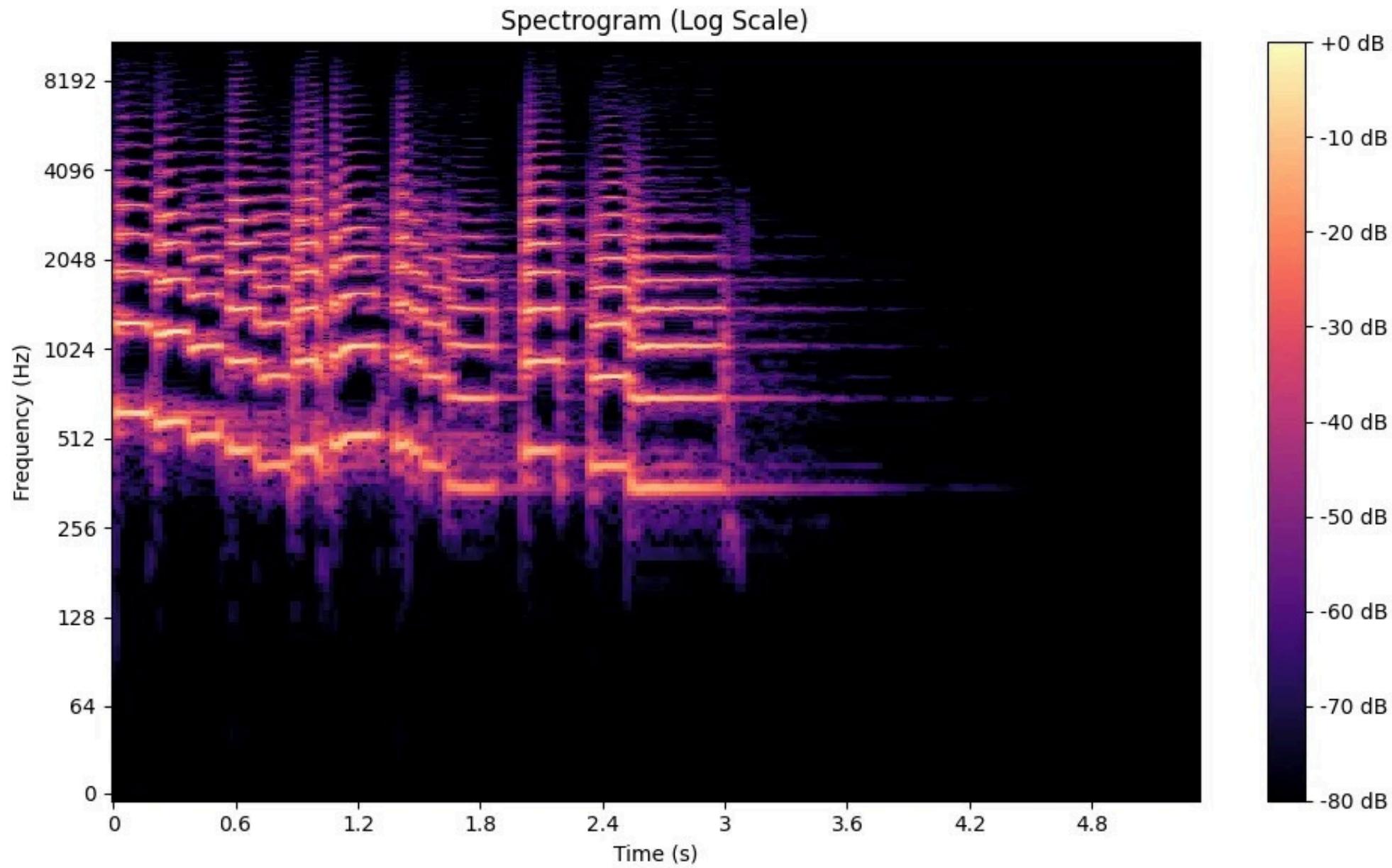
Types of Data

- Audio data



Types of Data

- Audio data



*

—

Data Literacy

Sources of Data

Sources of Data

kaggle

OpenML



Data Literacy

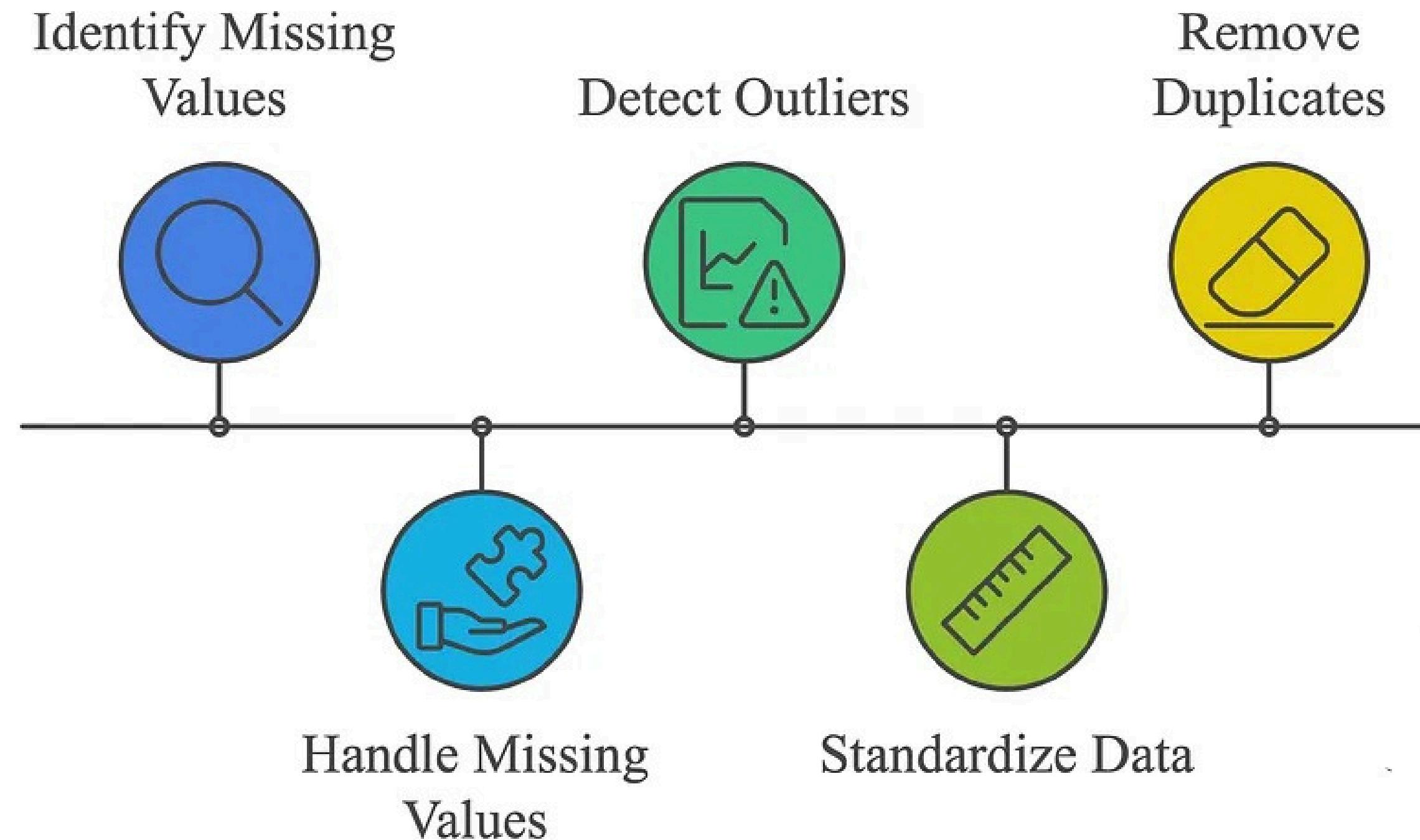
Data Cleaning

Data Cleaning

- **Data cleaning** is the process of fixing, correcting, and preparing raw data so it becomes accurate, consistent, and usable for analysis or machine learning.

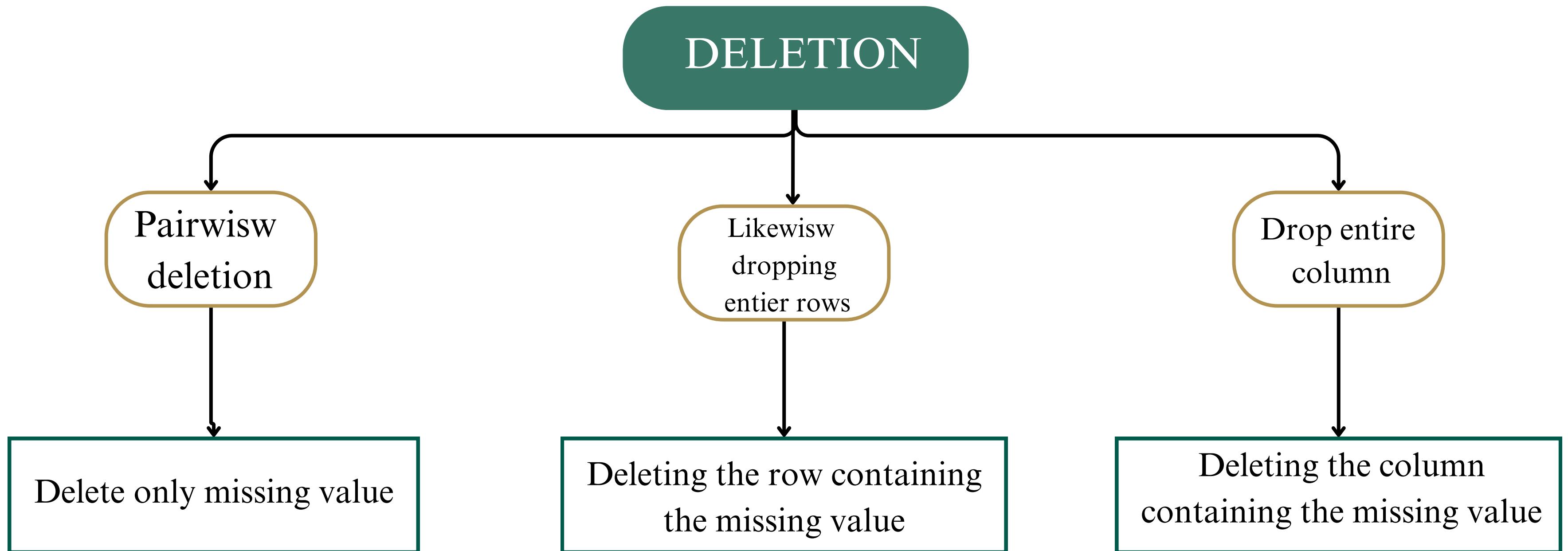


Data Cleaning



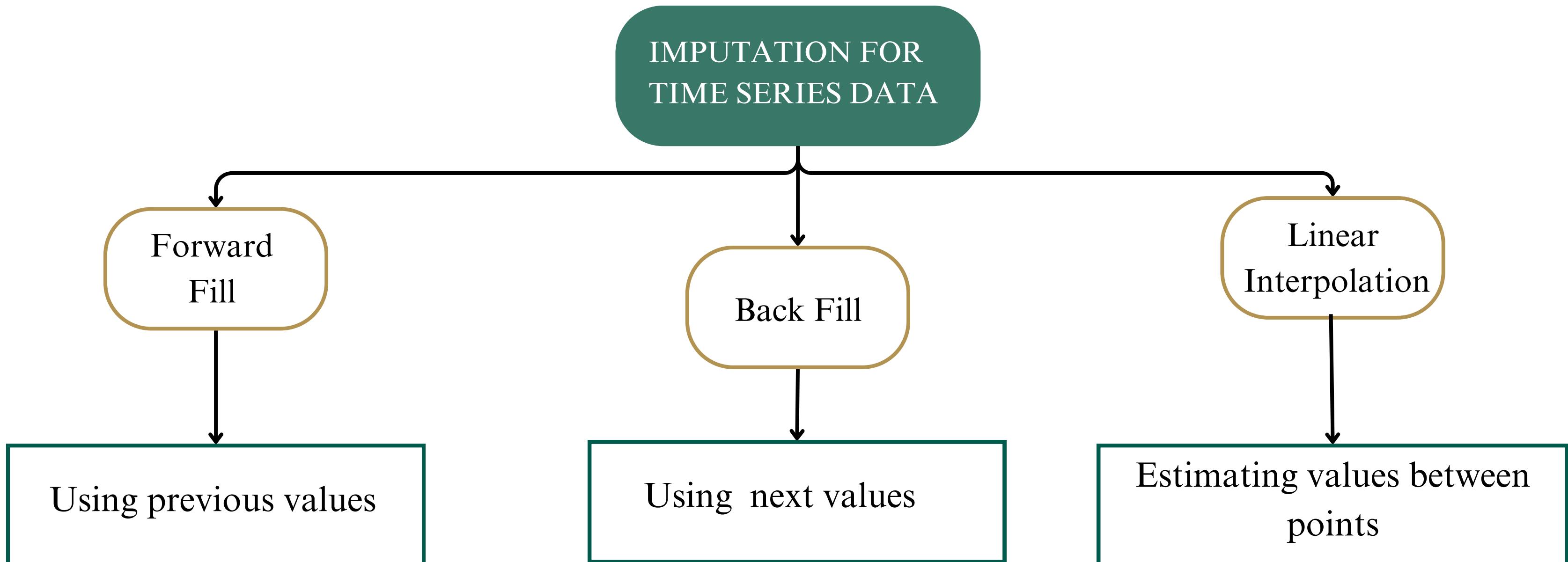
Data Cleaning

- Missing Values



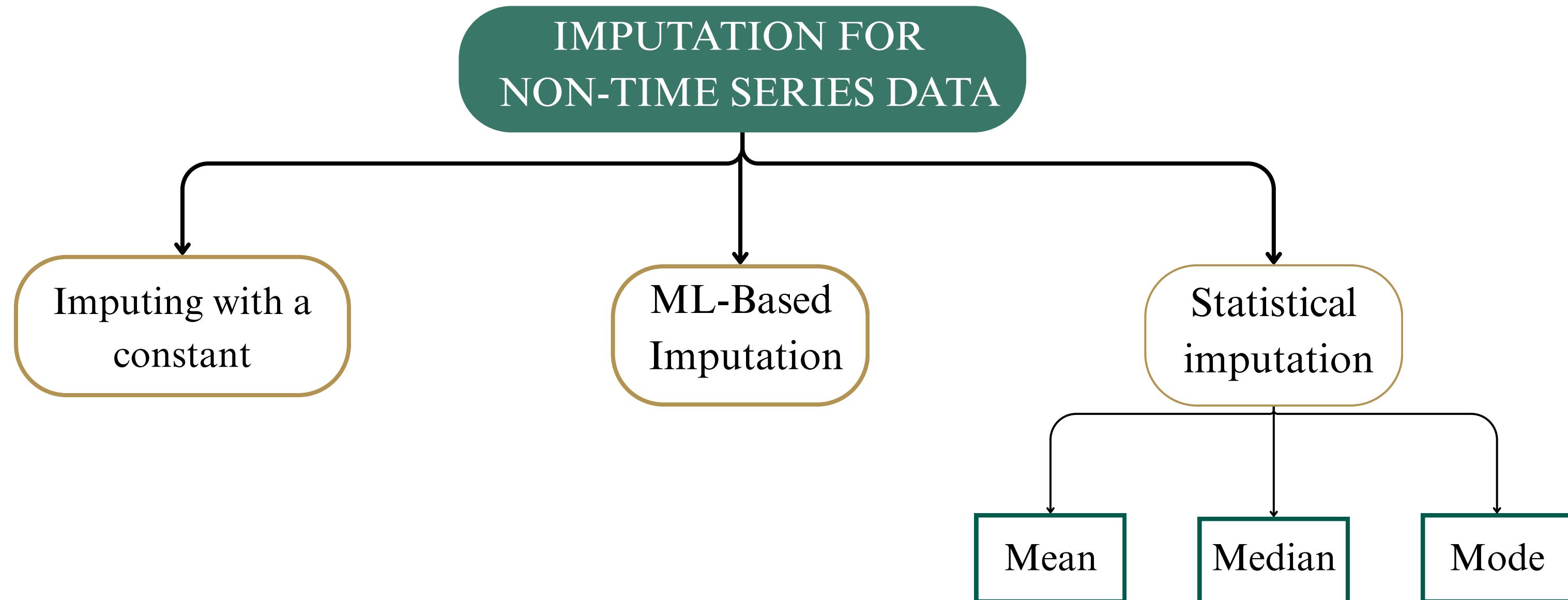
Data Cleaning

- Missing Values



Data Cleaning

- Missing Values

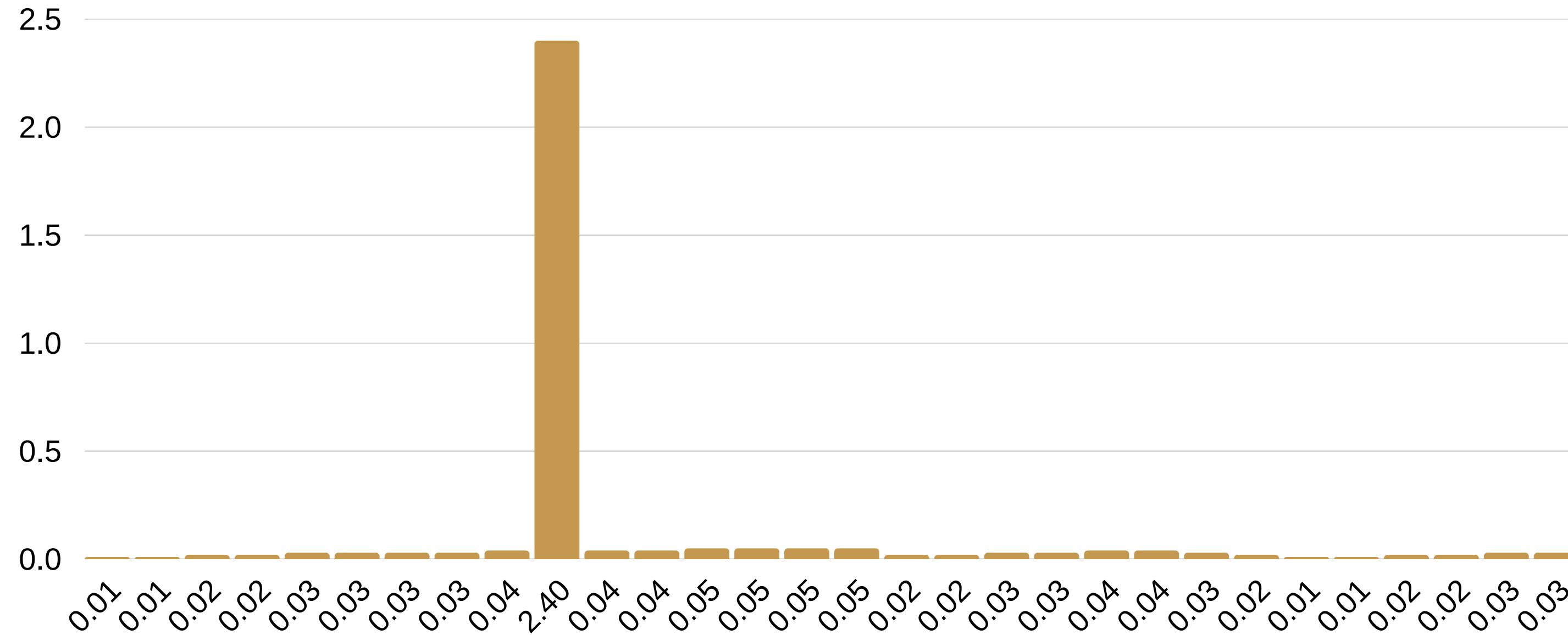


Does the type of
imputation matter



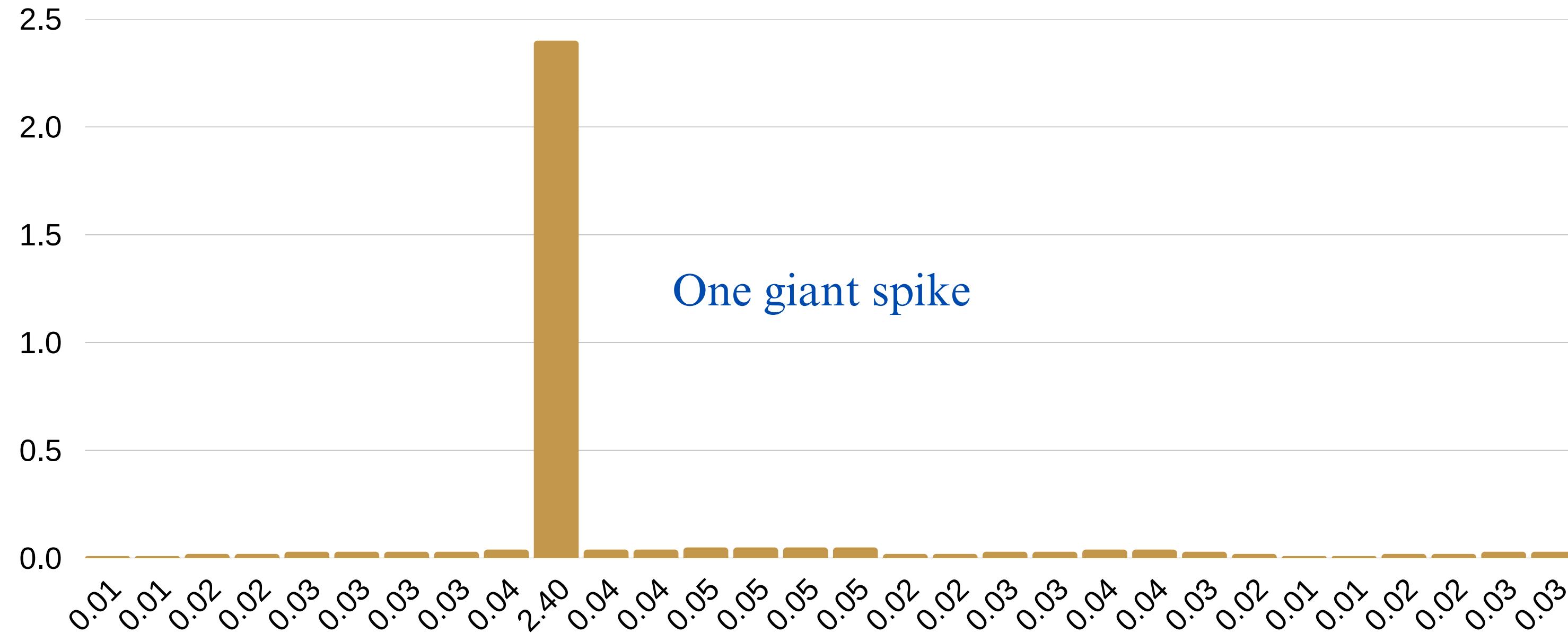
Data Cleaning

What's wrong with this histogram?



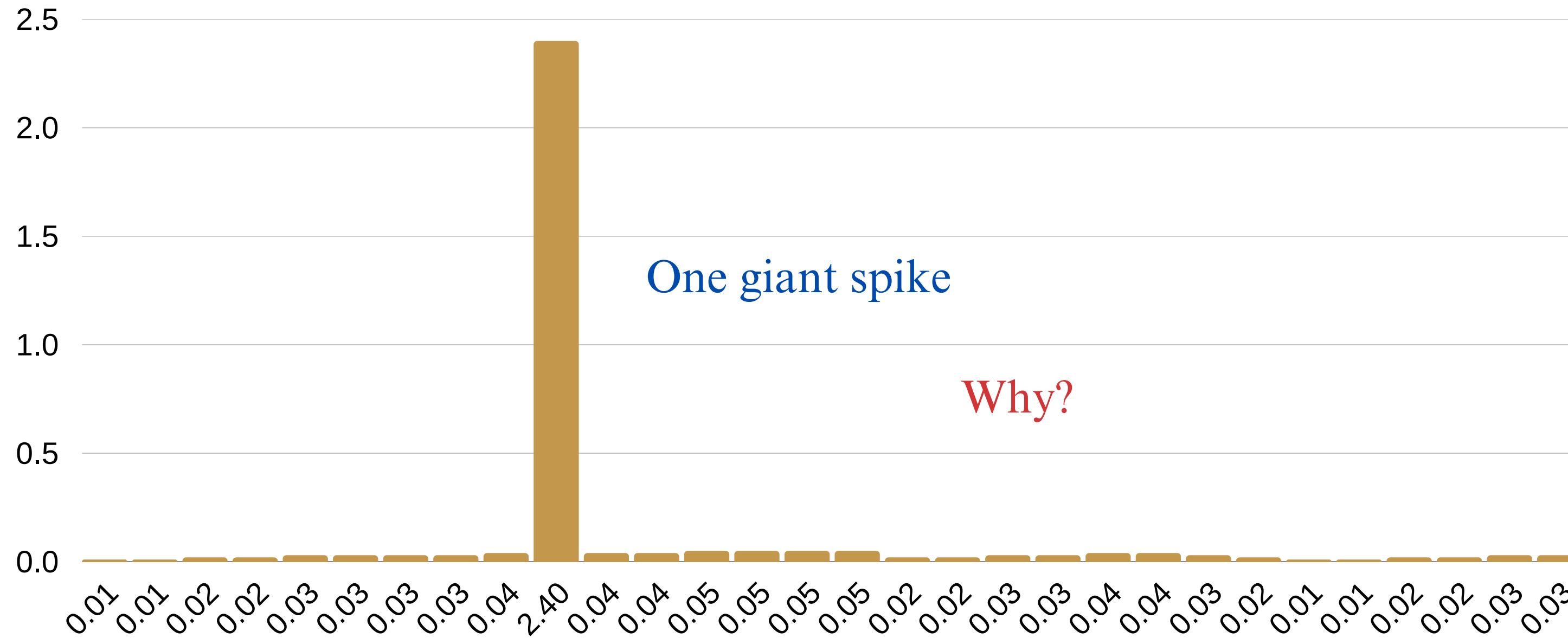
Data Cleaning

What's wrong with this histogram?



Data Cleaning

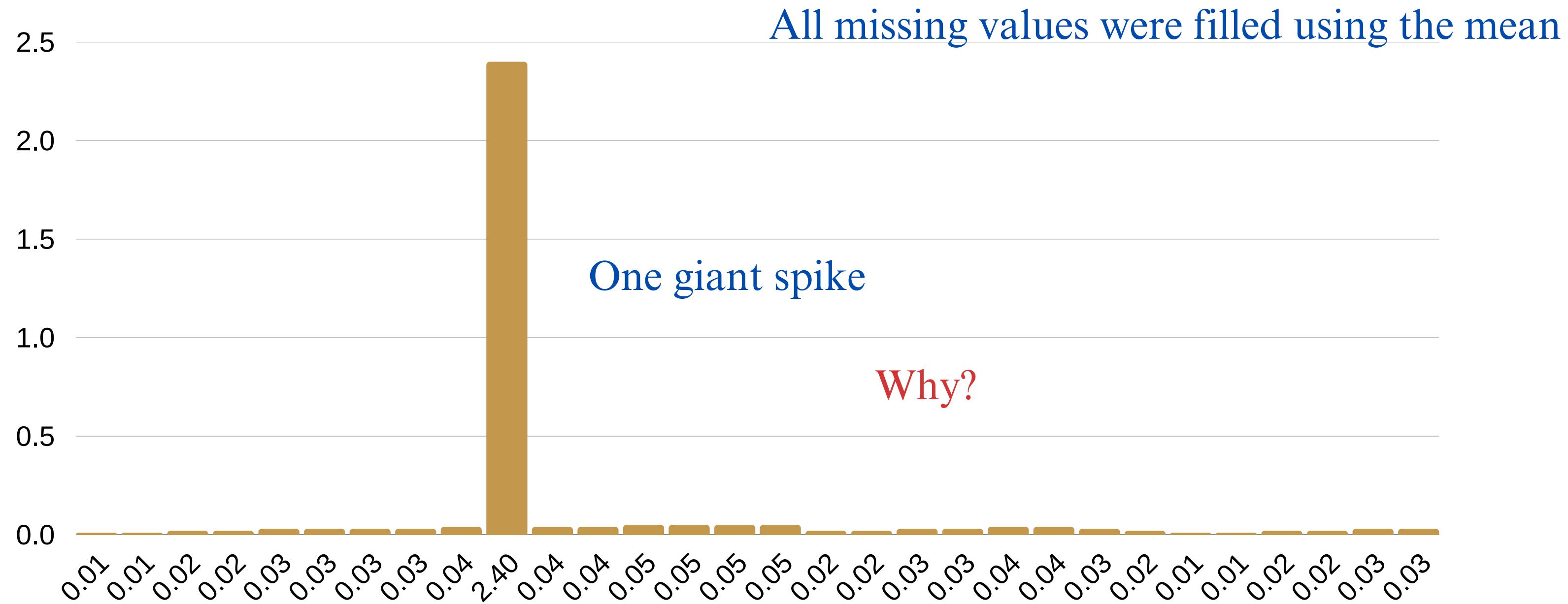
What's wrong with this histogram?



One giant spike
Why?

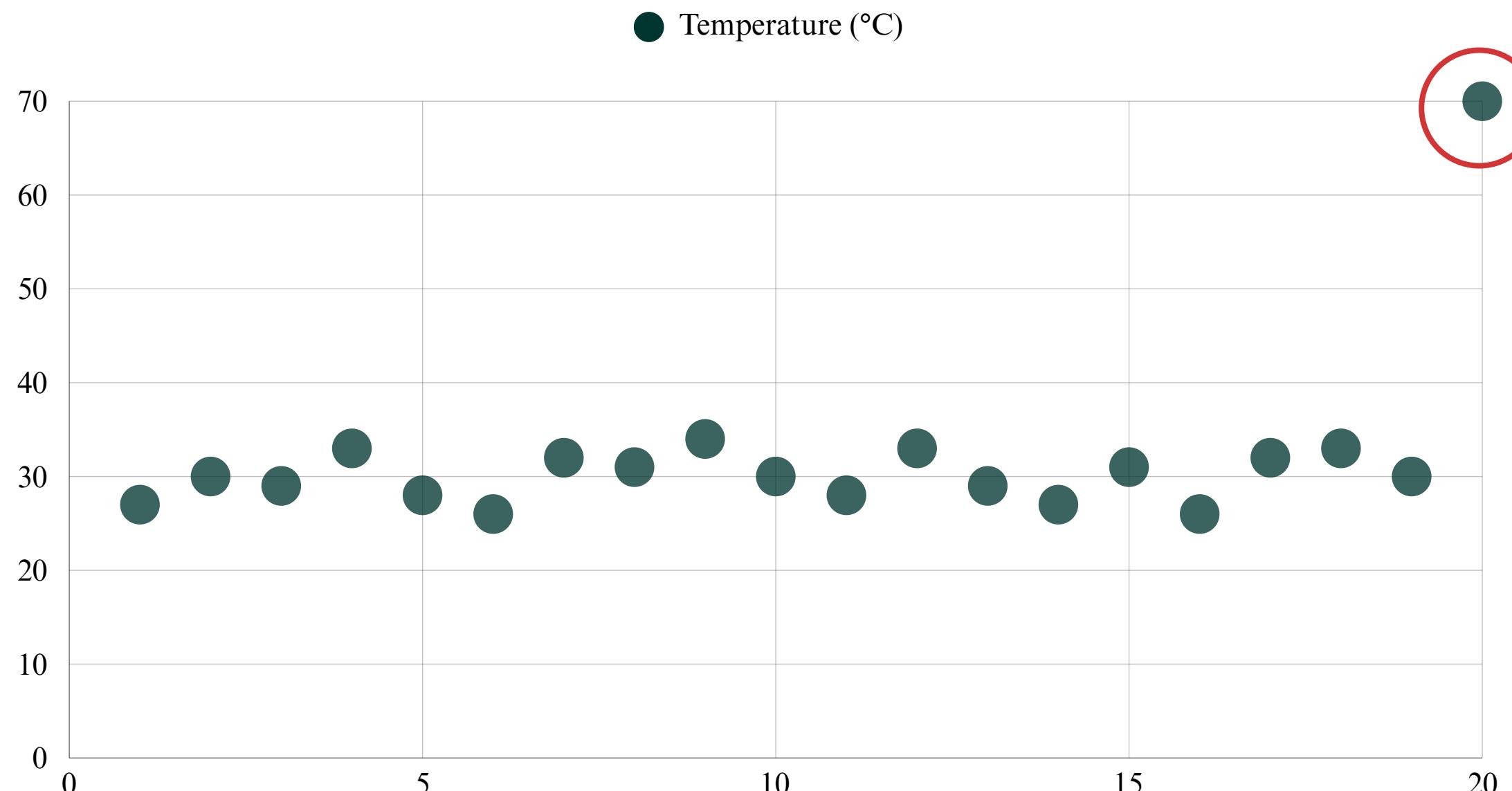
Data Cleaning

What's wrong with this histogram?



Data Cleaning

- Outliers

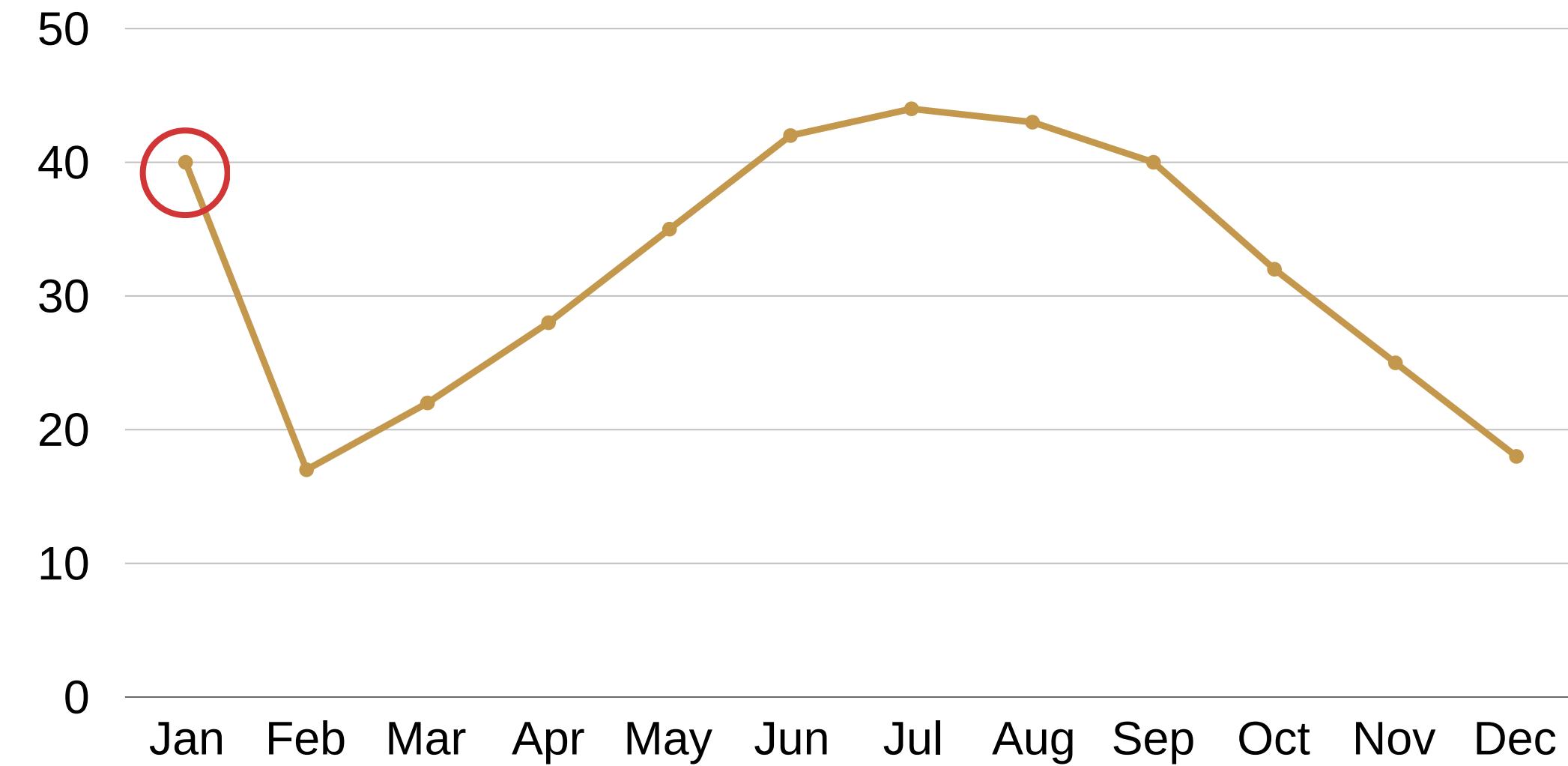


- Point Outliers

90% of temperatures are 25–35°C , one value is 70°C.

Data Cleaning

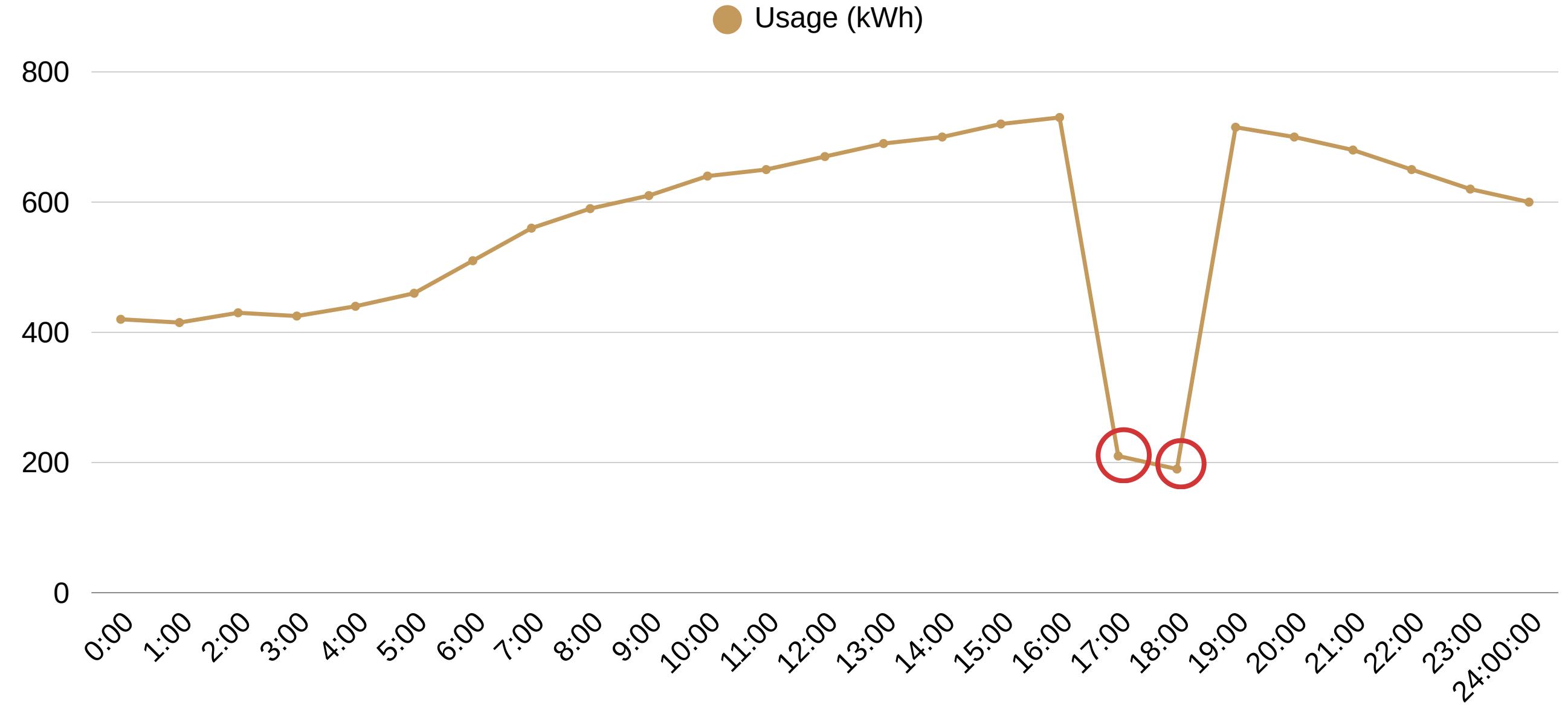
- Contextual Outliers



40°C normal in Riyadh summer, but 40°C in January becomes an outlier

Data Cleaning

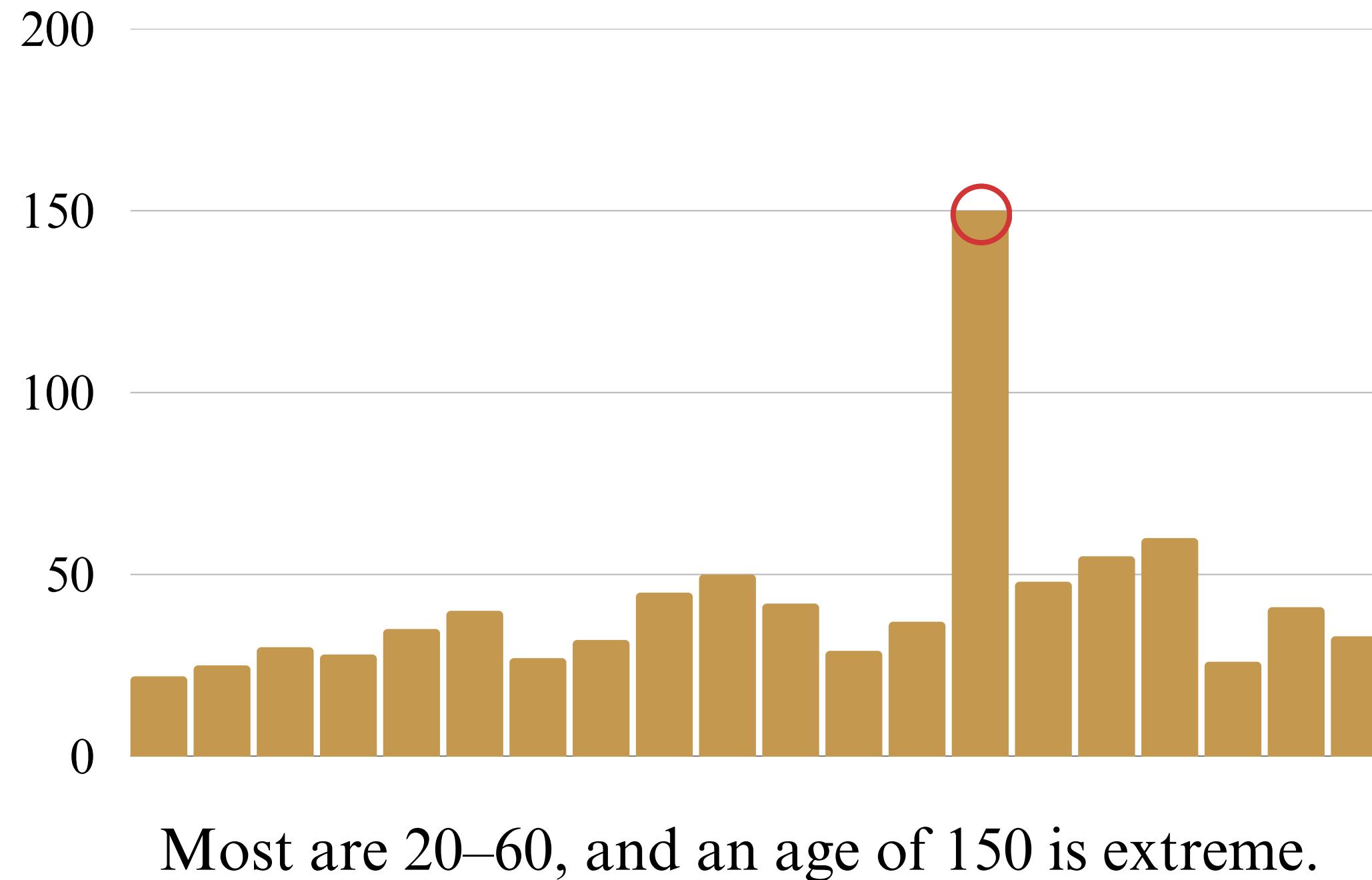
- Collective Outliers



Sudden drop in electricity usage for 2 hours, unusual as a chunk.

Data Cleaning

- Extreme Value Outliers



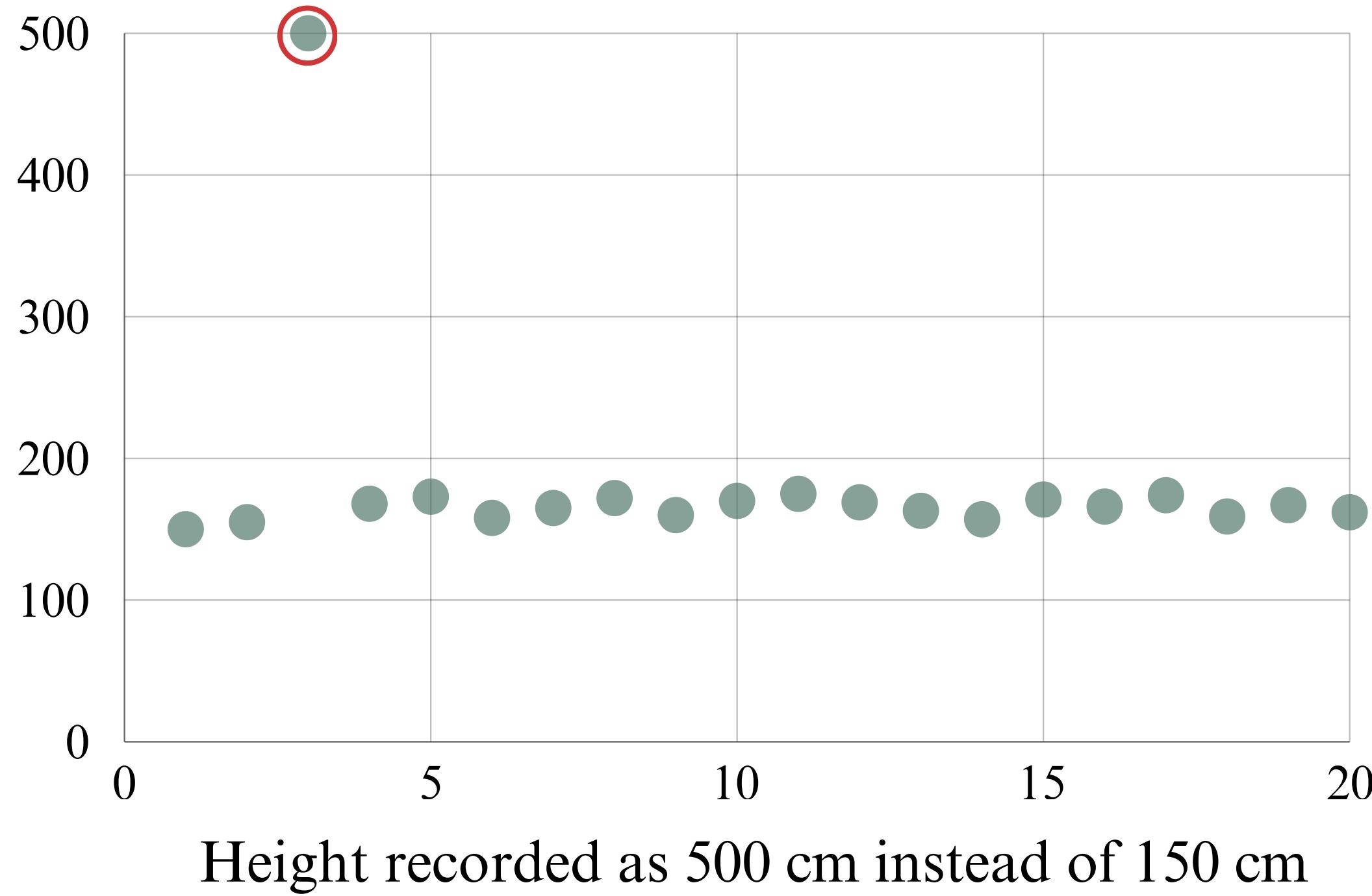
Data Cleaning

- Natural Outliers



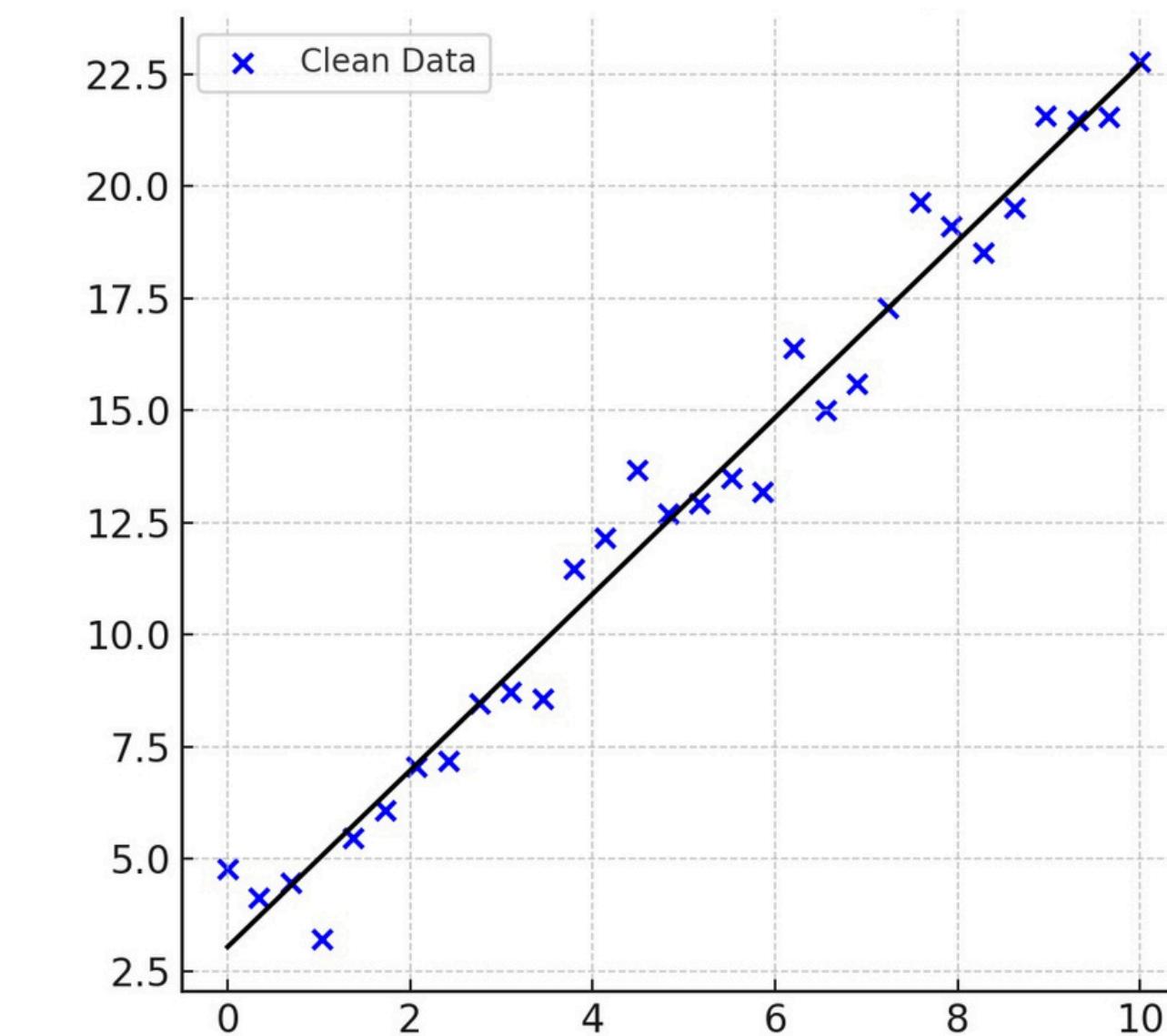
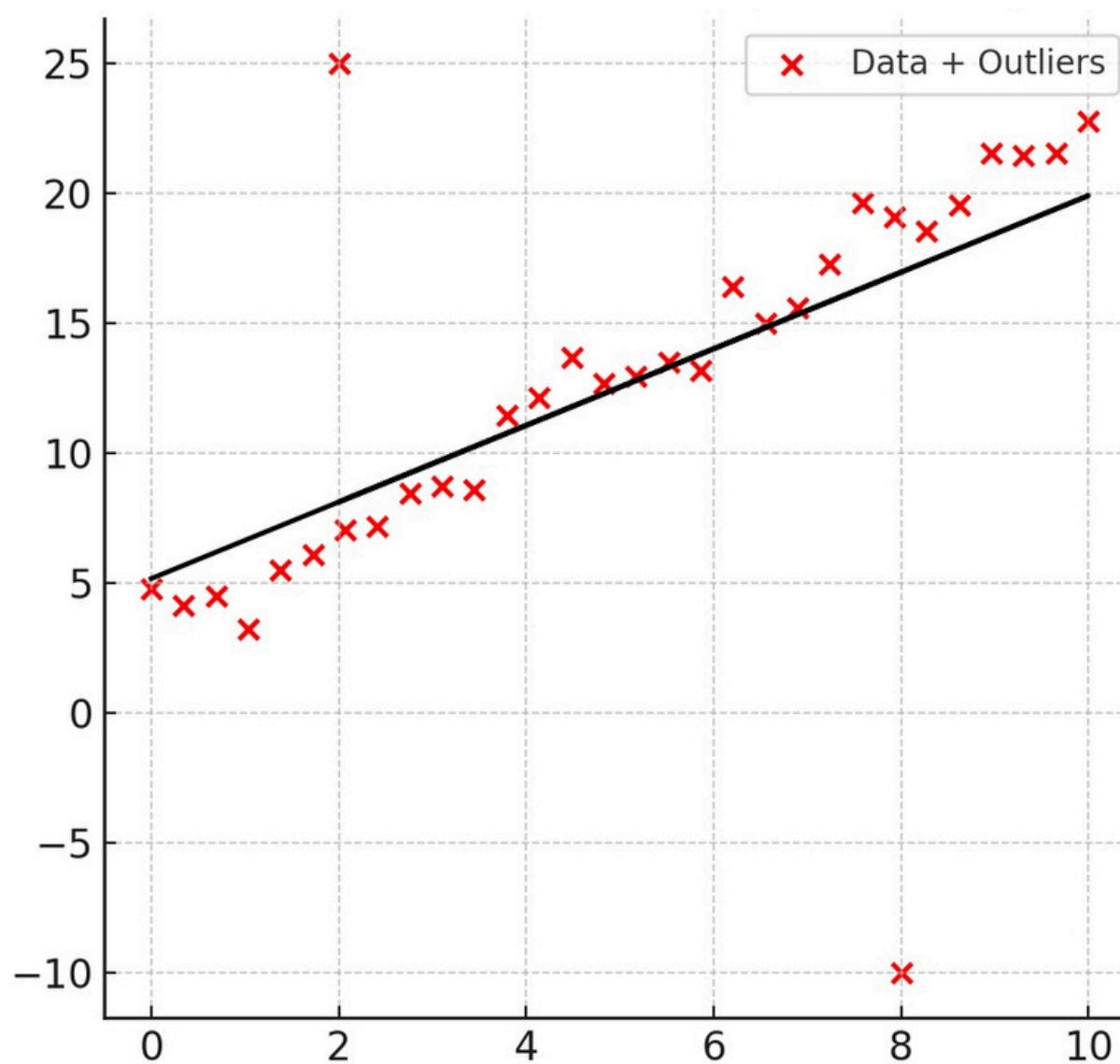
Data Cleaning

- Error Outlier



Data Cleaning

How outliers affect the behavior of a machine learning model.



Data Literacy

Hands-On Activity

Hands-On Activity

[Egypt Traffic Violations Dataset \(2024\).](#)



Hands-On Activity

HR Analytics



Materials

Day 2 Materials

Workshop Materials day 2



Quiz part 1

[quiz part 1](#)

quiz part 1



Quiz part 2

quiz part 2

quiz part 2





Thank You