

Lecture: Data Science: Pre-processing, Explorative Data Analysis

Data Science
Elf Data Science Courses, version 1.2.0

By Dr. Benjamin M. Henrich

Table of Content

- Pre-Alignment
- Theory: Pre-Processing - What is That?
- Case Study: Exploratory Data Analysis

Pre-alignment



After-Work-Drink? Wunschgetränk

Survey

efl | the Data Science Institute

Evaluierung - efl Data Science Courses

Wir würden uns sehr freuen, wenn Sie sich fünf bis zehn Minuten für die Evaluierung der Kurse Zeit nehmen würden.

[In Google anmelden](#), um den Fortschritt zu speichern. [Weitere Informationen](#)

* Gibt eine erforderliche Frage an

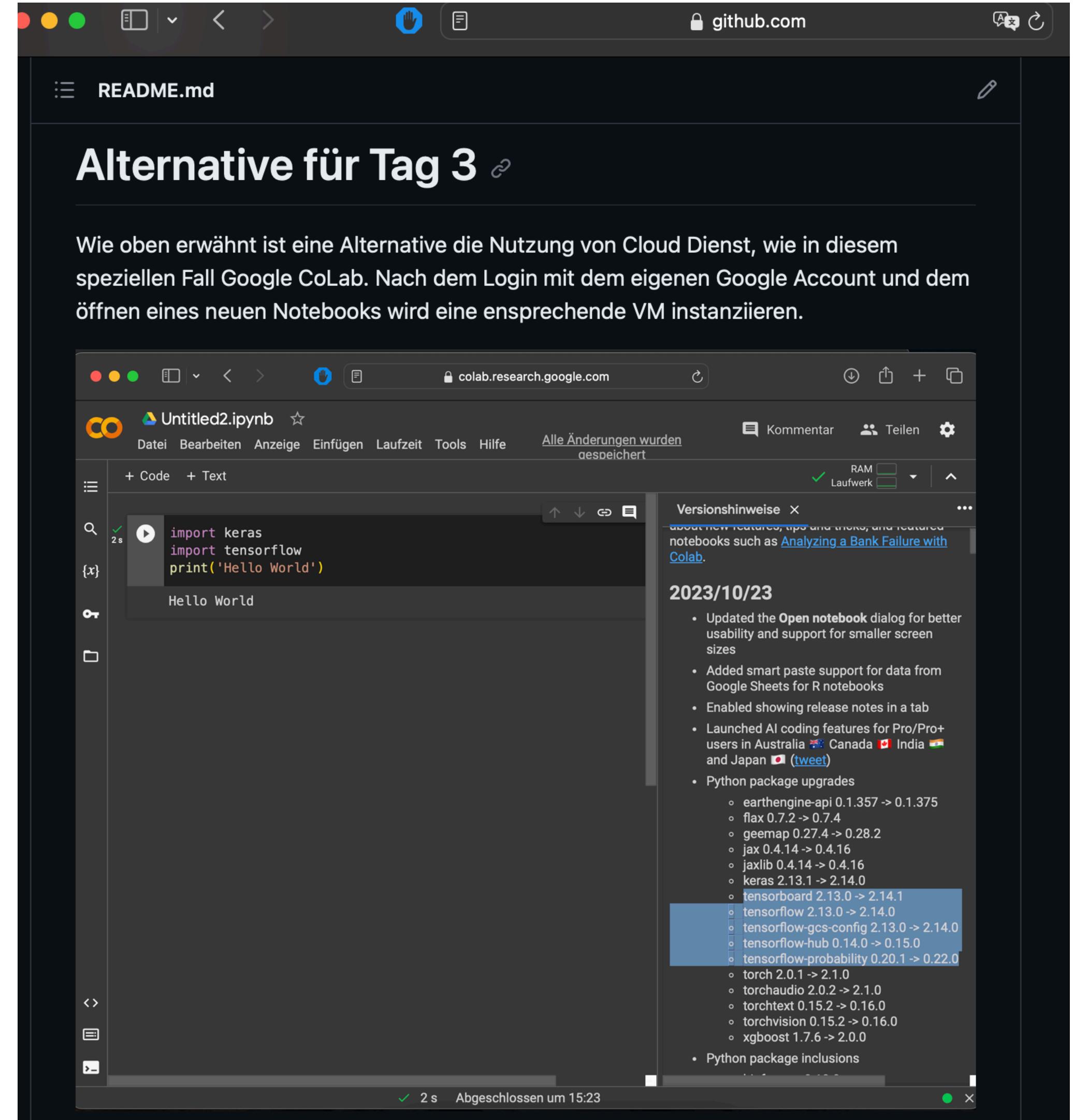
Gesamtbewertung



https://docs.google.com/forms/d/e/1FAIpQLSdH4RcSBN20rMX3bQjzfIt1JX_RjTijpiRebA_c062j83RFig/viewform

Alternative for Day 3

Google CoLab



github.com

README.md

Alternative für Tag 3

Wie oben erwähnt ist eine Alternative die Nutzung von Cloud Dienst, wie in diesem speziellen Fall Google CoLab. Nach dem Login mit dem eigenen Google Account und dem öffnen eines neuen Notebooks wird eine entsprechende VM instanzieren.

colab.research.google.com

Untitled2.ipynb

```
import keras
import tensorflow
print('Hello World')
```

2023/10/23

- Updated the Open notebook dialog for better usability and support for smaller screen sizes
- Added smart paste support for data from Google Sheets for R notebooks
- Enabled showing release notes in a tab
- Launched AI coding features for Pro/Pro+ users in Australia 🇦🇺 Canada 🇨🇦 India 🇮🇳 and Japan 🇯🇵 ([tweet](#))
- Python package upgrades
 - earthengine-api 0.1.357 -> 0.1.375
 - flax 0.7.2 -> 0.7.4
 - geemap 0.27.4 -> 0.28.2
 - jax 0.4.14 -> 0.4.16
 - jaxlib 0.4.14 -> 0.4.16
 - keras 2.13.1 -> 2.14.0
 - tensorboard 2.13.0 -> 2.14.1
 - tensorflow 2.13.0 -> 2.14.0
 - tensorflow-gcs-config 2.13.0 -> 2.14.0
 - tensorflow-hub 0.14.0 -> 0.15.0
 - tensorflow-probability 0.20.1 -> 0.22.0
 - torch 2.0.1 -> 2.1.0
 - torchaudio 2.0.2 -> 2.1.0
 - torchtext 0.15.2 -> 0.16.0
 - torchvision 0.15.2 -> 0.16.0
 - xgboost 1.7.6 -> 2.0.0
- Python package inclusions

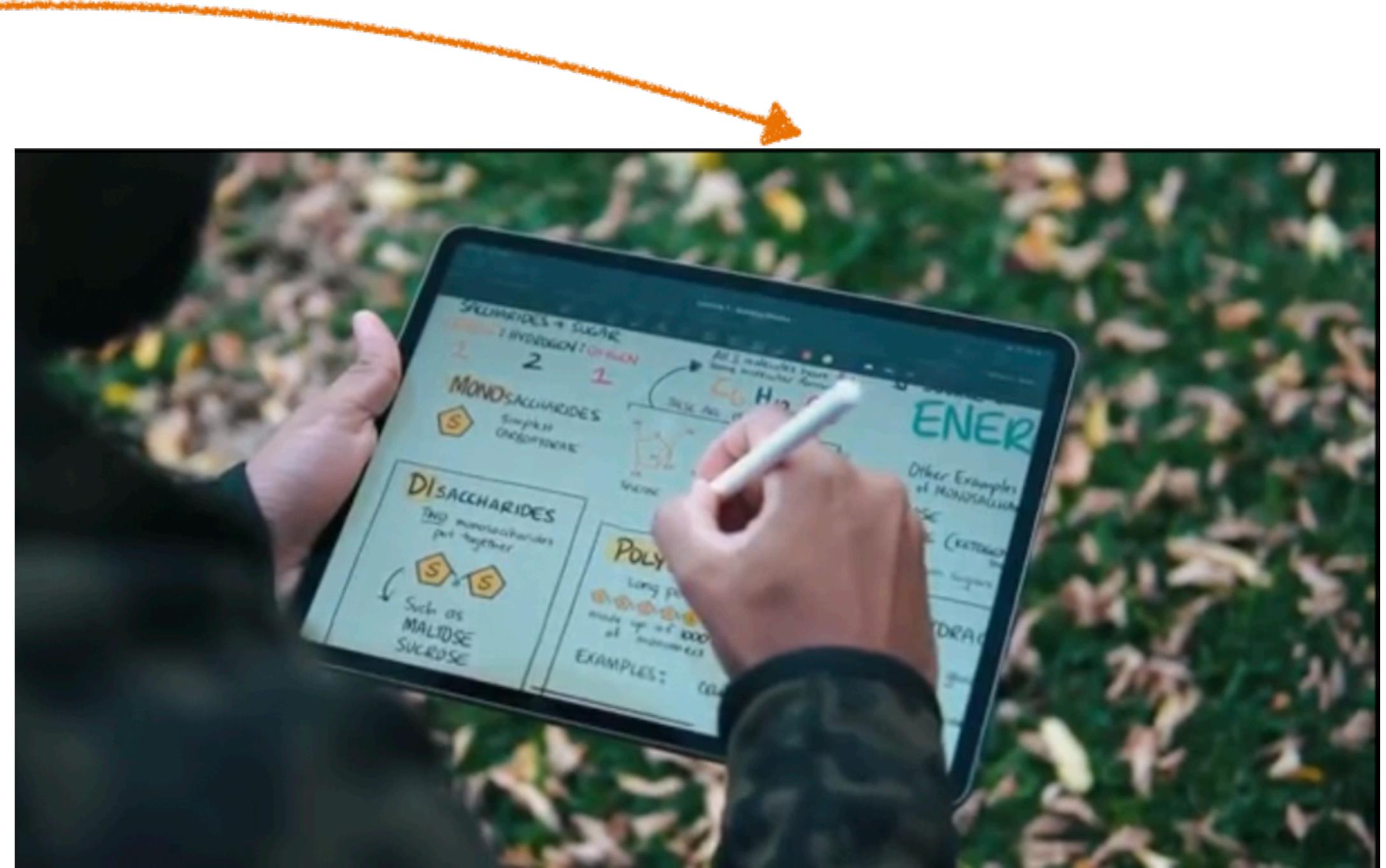
After-Work-Drink? Wunschgetränk



Table of Content

- Pre-Alignment
- Theory: Pre-Processing - What is That?
- Case Study: Exploratory Data Analysis

Data Science: But Why?

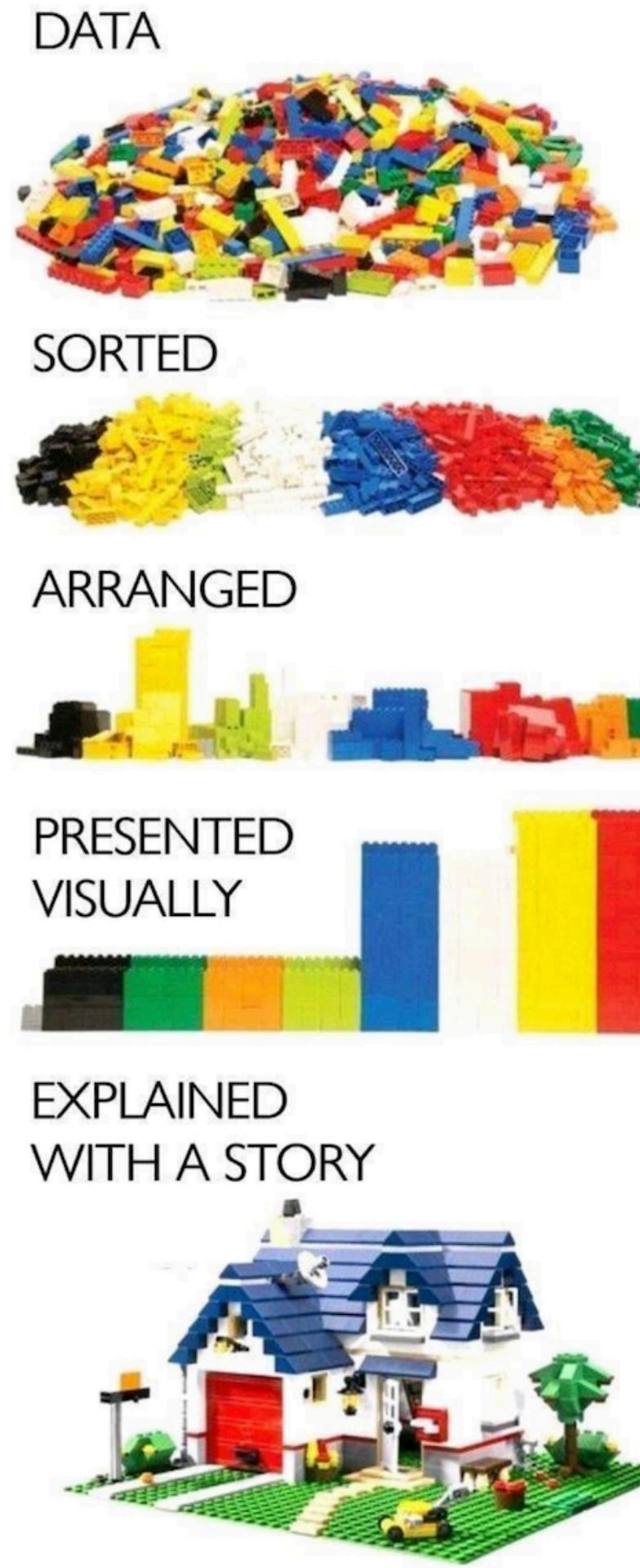


[1]

Raw Data

[2]

Knowlege



[1] <https://www.google.com/url?sa=i&url=https%3A%2F%2Flustich.de%2Fbilder%2Fandere%2Fpapierberg%2F&psig=AOvVaw1shfkEZVN9kWASoXICbwAi&ust=1601110332932000&source=images&cd=vfe&ved=2ahUKEwiMxrL29oPsAhVUNewKHdqPD0MQr4kDegUIARDIAQ>

Pre-Processing: What is That?

- Data pre-processing is a major and essential stage whose main goal is to **obtain final data sets** which can be considered correct and useful for further algorithms [1].
- The idea is **to filter the data** according to **quality criteria**, selected according to **possible target values**, and transformed correspondingly [2] .

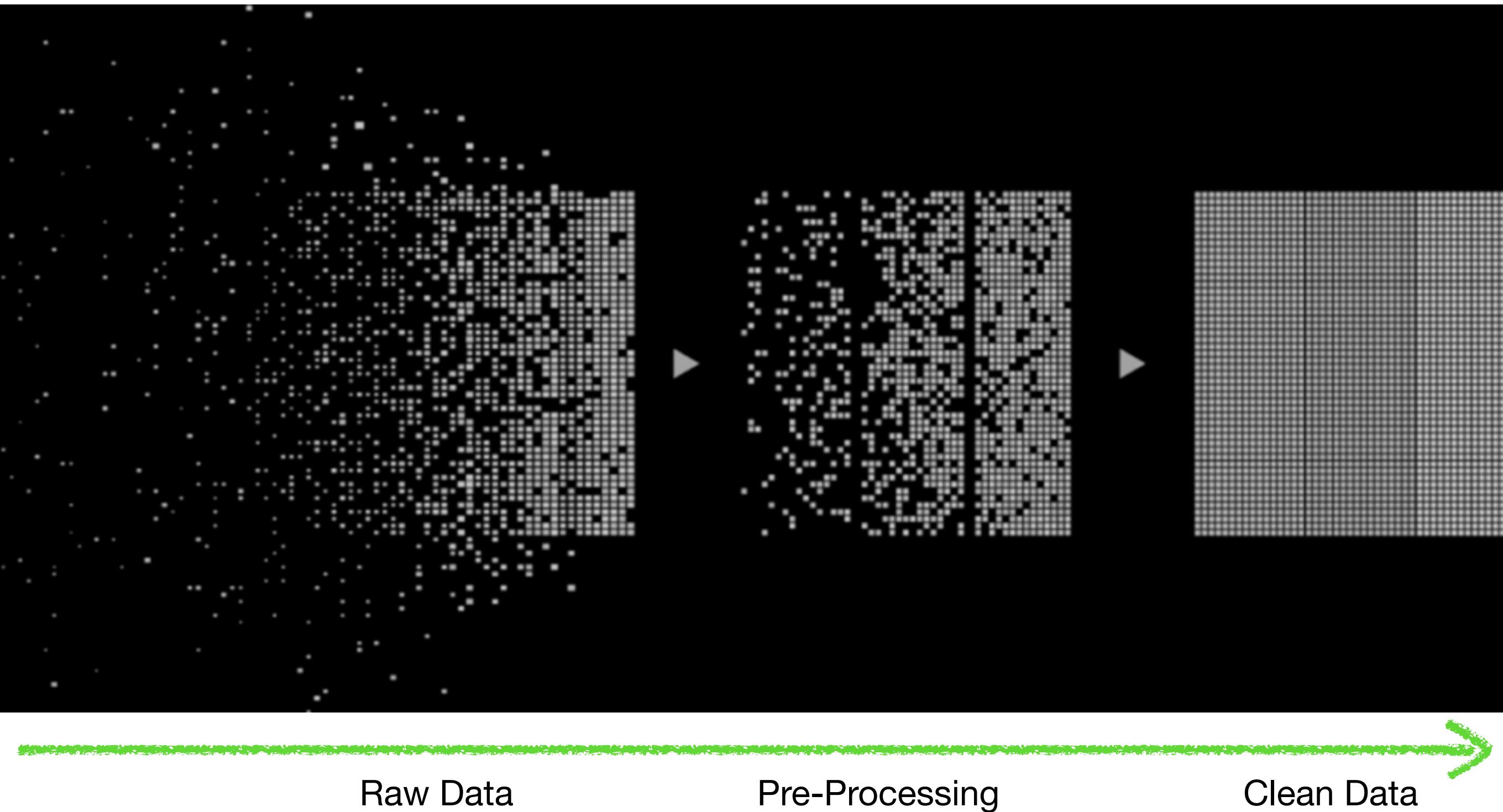
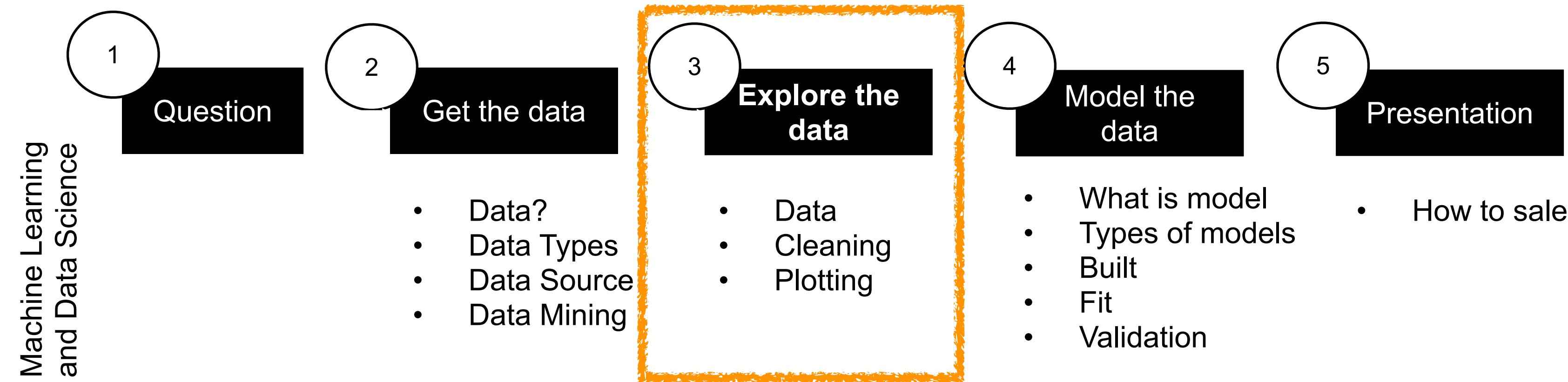


Figure: From Raw Data to Clean Data

[1] García, S., Ramírez-Gallego, S., Luengo, J., Benítez, J. M., & Herrera, F. (2016). Big data preprocessing: methods and prospects. *Big Data Analytics*, 1(1), 1-22.

[2] Abdel-Karim, B. M., Pfeuffer, N., & Hinz, O. (2021). Machine learning in information systems-a bibliographic review and open research issues. *Electronic Markets*, 31, 643-670.

Pre-Processing as Part of the KDD Model



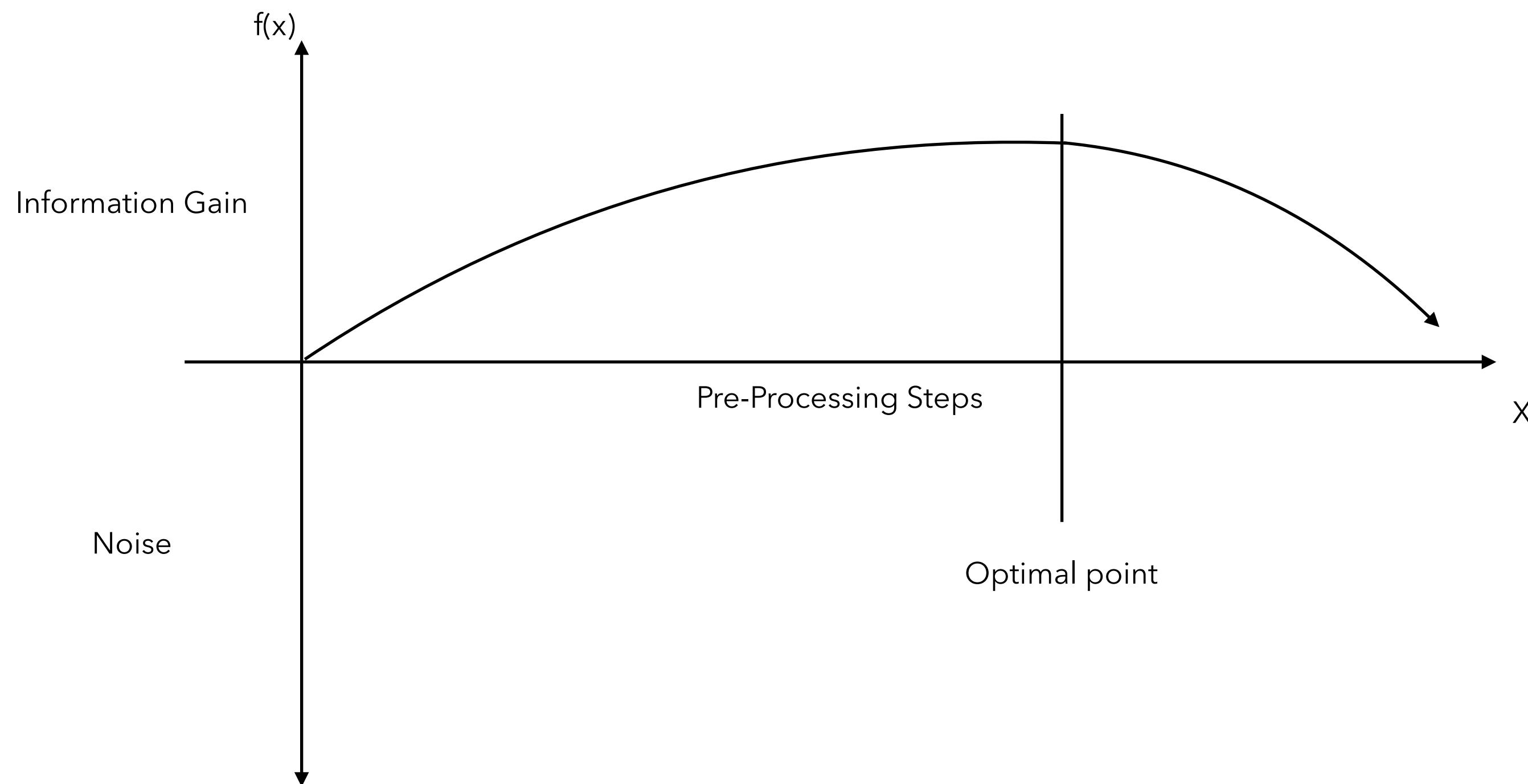
KDD Model: Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). From data mining to knowledge discovery in databases. *AI magazine*, 17(3), 37-37.

Pre-Processing: General Node

- The preprocessing steps are highly dependent on the modeling approaches used.
- **There is no right or best preprocessing for all approaches**
- For example:
 - What measurement scales are excepted?
 - Can the approach deal with missing values
 - Is the approach robust with regard to outliers?
 - Is the approach robust with regard to differences in ranges?
 - Is the approach robust with regard to irrelevant features?
 - Is the approach robust with regard to redundant features?

Mishra, P., Biancolillo, A., Roger, J. M., Marini, F., & Rutledge, D. N. (2020). New data preprocessing trends based on ensemble of multiple preprocessing techniques. *TrAC Trends in Analytical Chemistry*, 132, 116045.

Data Pre-Processing is a Trade-Off: Between Increasing Information Gain and Introducing Noise



Chicco, D. (2017). Ten quick tips for machine learning in computational biology. *BioData Mining*, 10(1), 35
Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer. [Link](#)

Overview

- **Handling Missing Data**
 - Removing missing values
 - Filling missing values
- **Data Type Conversion**
- **Handling Duplicates**
- **Encoding Categorical Variables**
 - Label Encoding
 - One-Hot Encoding
 - ...
- Feature Scaling
 - Normalization
 - Standardization
- Outlier Detection and Removal

Alasadi, S. A., & Bhaya, W. S. (2017). Review of data preprocessing techniques in data mining. *Journal of Engineering and Applied Sciences*, 12(16), 4102-4107.

Fan, C., Chen, M., Wang, X., Wang, J., & Huang, B. (2021). A review on data preprocessing techniques toward efficient and reliable knowledge discovery from building operational data. *Frontiers in energy research*, 9, 652801.

García, S., Ramírez-Gallego, S., Luengo, J., Benítez, J. M., & Herrera, F. (2016). Big data preprocessing: methods and prospects. *Big data analytics*, 1, 1-22.

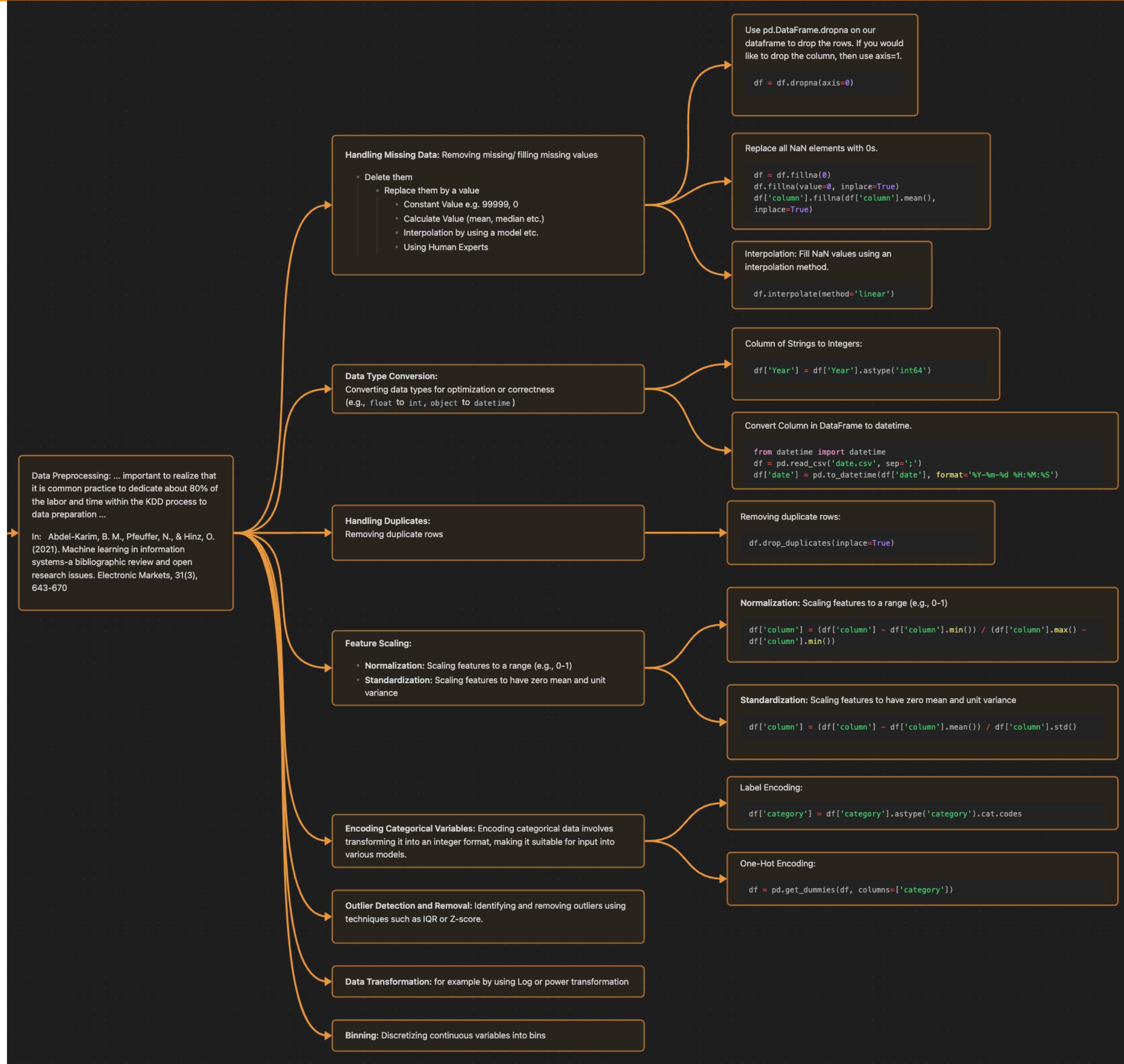
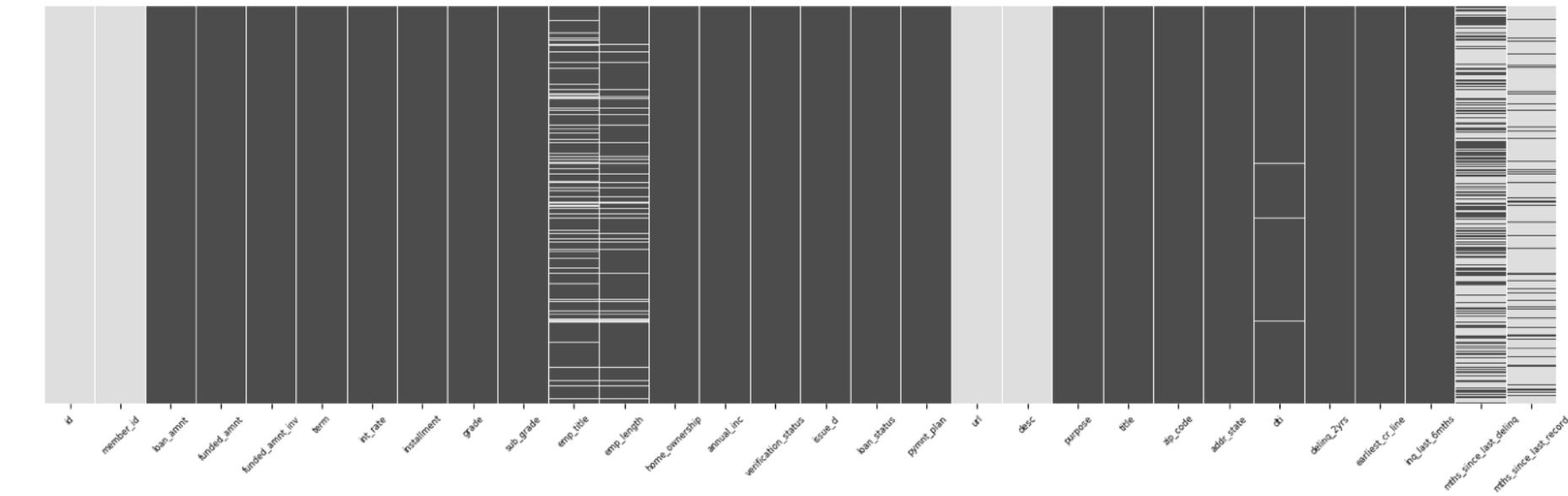


Figure: Obsidian Mindmap

Handling Missing Data

- Find all missing values.
- We should get an overview of the missing values.

```
plt.figure(figsize=(15, 5))
sns.heatmap(df.isnull(), cbar=False,
yticklabels=False, cmap='winter')
plt.xticks(rotation=45, fontsize=6)
plt.tight_layout()
plt.savefig('figures/
MissingValues.png')
plt.show()
```



Important Notes about NA/ NAN/ Not-a-Number in Python in regards to the data types:

- In Python, you can create nan with float('nan'), math.nan, or np.nan. nan is considered a missing value in pandas., s_nan = pd.Series([float('nan')], math.nan, np.nan])
- None is also considered a missing value
- nan is a floating-point number float
- None in the object column remains as None
- String is not considered a missing value
- the strings 'NaN' and 'None' are not treated as missing values. The empty string "" is also not considered a missing value.
- if you want to treat certain values as missing, you can use the replace() method to replace them with float('nan'), np.nan, or math.nan., such as s_replace = s_str.replace(['NaN', 'None', ''], float('nan'))
s_replace = s_str.replace(['NaN', 'None', ''], float(np.nan))

Handling Missing Data: Removing

```
import pandas as pd
import numpy as np

# Sample DataFrame with missing values
data = {
    'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
    'Age': [55, np.nan, 35, 45, np.nan],
    'City': ['London', 'Los Angeles', np.nan, 'Chicago', 'New York']
}

df = pd.DataFrame(data)
print(df)
```

	Name	Age	City
0	Alice	55.0	London
1	Bob	NaN	Los Angeles
2	Charlie	35.0	NaN
3	David	45.0	Chicago
4	Eve	NaN	New York

Remove rows with any missing values:

df_dropped_rows = df.dropna() Deleting:

Remove columns with any missing values:

df_dropped_columns = df.dropna(axis=1) Deleting:



Figure: Meme About Replacing NaN

Handling Missing Value

```
import pandas as pd
import numpy as np

# Sample DataFrame with missing values
data = {
    'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
    'Age': [55, np.nan, 35, 45, np.nan],
    'City': ['London', 'Los Angeles', np.nan, 'Chicago', 'New York']
}

df = pd.DataFrame(data)
print(df)
```

	Name	Age	City
0	Alice	55.0	London
1	Bob	NaN	Los Angeles
2	Charlie	35.0	NaN
3	David	45.0	Chicago
4	Eve	NaN	New York

Fill with a specific value:

```
df_filled = df.fillna(0)
```

Forward fill (propagate last valid observation forward):

```
df_filled_forward = df.fillna(method='ffill')
```

Backward fill (propagate next valid observation backward):

```
df_filled_backward = df.fillna(method='bfill')
```

Note: Interpolate missing data, which is especially useful for time series.

NaN => 0

Bob, **55.0**, Los Angeles
 Charlie, 35.0, **Los Angeles**
 Eve, **45.0**, New York

Careful

Bob, **35.0**, Los Angeles
 Charlie, 35.0, **Chicago**
 Eve, **NaN**, New York



Handling Missing Values: Overview

Method	Description	Code Example
Check for Missing Values	Identify missing values in the DataFrame	<code>df.isnull()</code>
Count Missing Values	Count missing values in each column	<code>df.isnull().sum()</code>
Drop Rows with NaN	Remove rows that have at least one missing value	<code>df.dropna()</code>
Drop Columns with NaN	Remove columns with missing values	<code>df.dropna(axis=1)</code>
Fill Missing Values	Replace NaN with a specific value (e.g., 0, mean, median)	<code>df.fillna(value)</code>
Forward Fill	Propagate last valid value forward	<code>df.fillna(method='ffill')</code>
Backward Fill	Use next valid value to fill missing values	<code>df.fillna(method='bfill')</code>
Fill with Column Mean	Replace NaN with column mean	<code>df['col'].fillna(df['col'].mean())</code>
Interpolate	Estimate missing values using interpolation	<code>df.interpolate()</code>
Check for Missing Rows	Filter rows that contain missing values	<code>df[df.isnull().any(axis=1)]</code>
Replace NaN	Replace NaN with a specific value	<code>df.replace(np.nan, 'value')</code>
Custom Imputation	Use custom logic for filling missing values	<code>df['col'].fillna(custom_function())</code>
Count Non-Missing Values	Get count of non-null values in each column	<code>df.notnull().sum()</code>
Detect Missing Values in Series	Check for NaN in a specific column or Series	<code>df['col'].isnull()</code>
Drop Rows with All NaN	Drop rows where all elements are NaN	<code>df.dropna(how='all')</code>
Drop Rows with Threshold	Keep rows with a minimum number of non-NaN values	<code>df.dropna(thresh=2)</code>

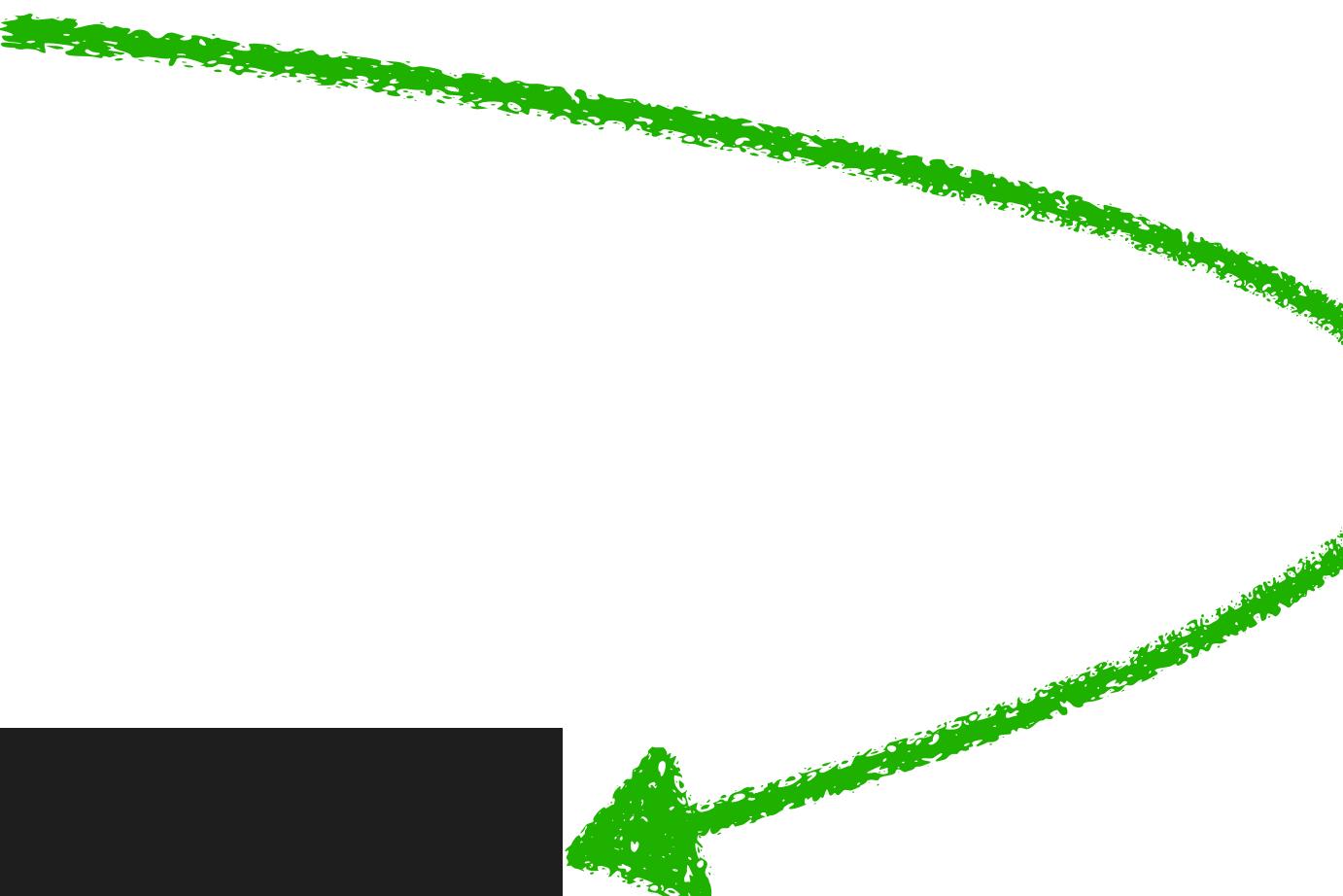
Notes:

- `fillna()` allows filling NaN values with constants, mean, median, or previous/next values.
- `dropna()` removes rows or columns containing NaN.
- `interpolate()` can estimate values using methods like linear, polynomial, or spline interpolation.

Data Type Conversion (Data Types of Each Column)

```
#   Column    Non-Null Count  Dtype  
---  --  
0   ID        4 non-null      object 
1   Price     4 non-null      object 
2   Date      4 non-null      object 
3   Category  4 non-null      int64  
dtypes: int64(1), object(3)  
memory usage: 260.0+ bytes
```

```
#   Column    Non-Null Count  Dtype  
---  --  
0   ID        4 non-null      int64  
1   Price     4 non-null      float64 
2   Date      4 non-null      datetime64[ns] 
3   Category  4 non-null      category
dtypes: category(1), datetime64[ns](1), float64(1), int64(1)  
memory usage: 356.0 bytes
None
```



```
data = {  
    'ID': ['1', '2', '3', '4'], # IDs stored as strings
    'Price': ['99.5', '20.99', '15.75', '40'], # Prices stored as strings
    'Date': ['2023-01-01', '2023-01-02', '2023-01-03', '2023-01-04'], # Dates stored as strings
    'Category': [1, 2, 1, 2] # Categories stored as integers
}

df = pd.DataFrame(data)
print("Original DataFrame:\n", df)

print(df.info())

# Convert ID from string to integer:
df['ID'] = df['ID'].astype(int)
print("\nDataFrame after converting ID to integer:\n", df)

# Convert Price from string to float:
df['Price'] = df['Price'].astype(float)
print("\nDataFrame after converting Price to float:\n", df)

# Convert Date from string to datetime:
df['Date'] = pd.to_datetime(df['Date'])
print("\nDataFrame after converting Date to datetime:\n", df)

# Convert Category from integer to category:
df['Category'] = df['Category'].astype('category')
print("\nDataFrame after converting Category to categorical:\n", df)

print(df.info())
```

Figure: Data Type Conversion

Data Type Conversion: But Why?

Correct Data Representation:

- Data types ensure that each column in a DataFrame is stored in the most appropriate format. For example:
- Numerical Data: Converting columns to int or float ensures that numerical operations (like addition or averaging) are performed correctly.
- Categorical Data: Converting columns to category type can save memory and speed up operations if the column has a limited number of unique values.
- Date/Time Data: Converting columns to datetime type allows for time-based operations and manipulations.

Performance Optimization:

- Correct data types can lead to more efficient storage and faster computation. For instance:
- Memory Efficiency: Using category for columns with a limited number of unique values reduces memory usage compared to using object (which is typically a catch-all for string data).
- Speed: Numerical operations are faster with int or float types compared to object types, which are more general-purpose and less optimized for numeric computations.

Handling Duplicates

Why handling duplicates could be important:

- **Data Integrity:** Duplicates can compromise the accuracy and integrity of the data.
- **Accurate Analysis:** impact on statistical analysis and metrics.
- **Efficient Ressourcen Usage**

Option to Handle:

- **Removing Duplicates:** Remove duplicate rows based on all columns or specific ones.
- **Identifying Duplicates:** Find duplicate rows before deciding how to handle them.



Figure: The Matrix: Agent Smith's Replicating

<https://www.cbr.com/the-matrix-agent-smith-replicating-powers-mysteries-explained/>

Label Encoding

Label Encoding is a technique used to convert categorical variables into numerical values. Each unique category in a categorical variable is assigned an integer value.

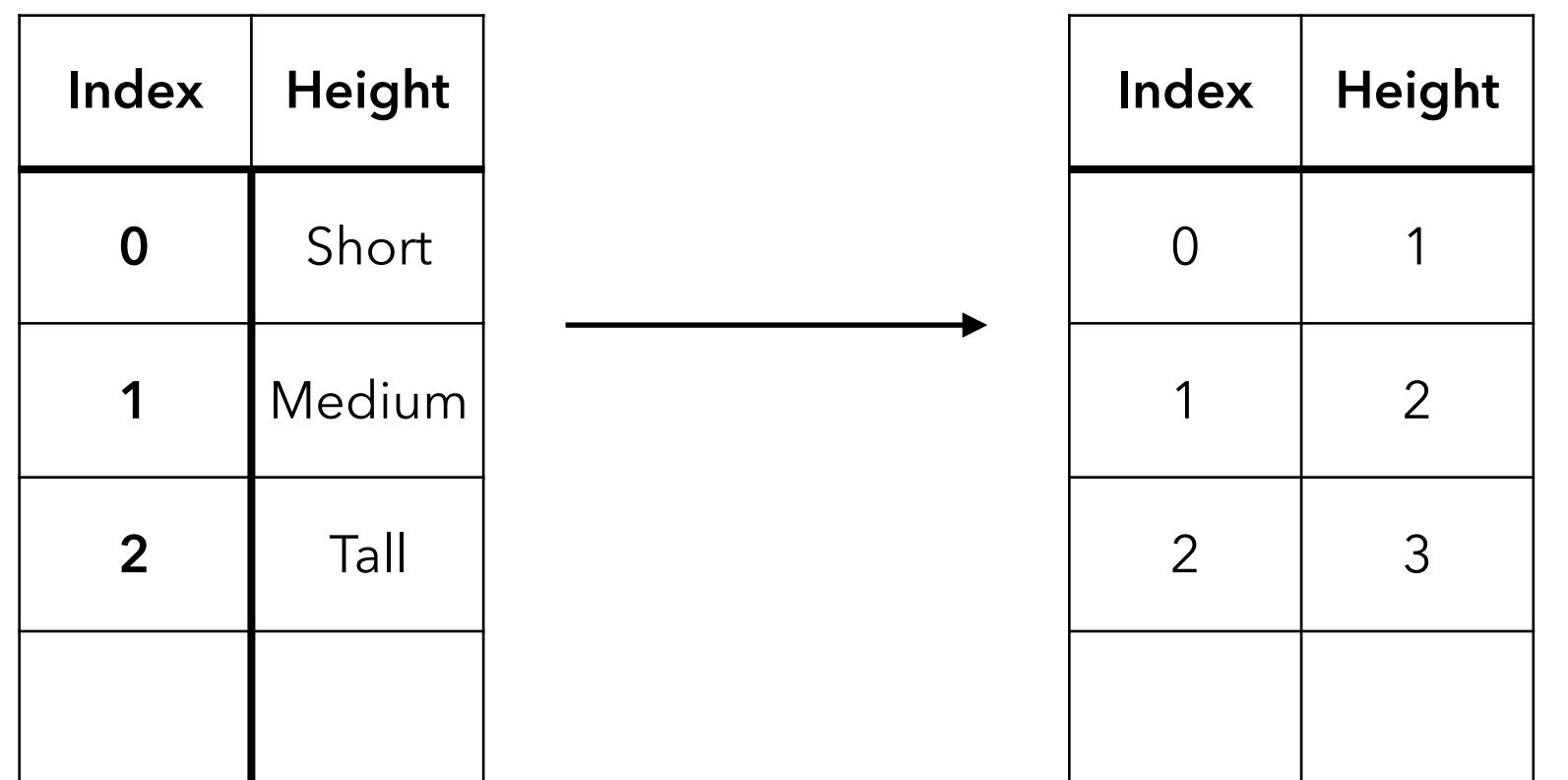
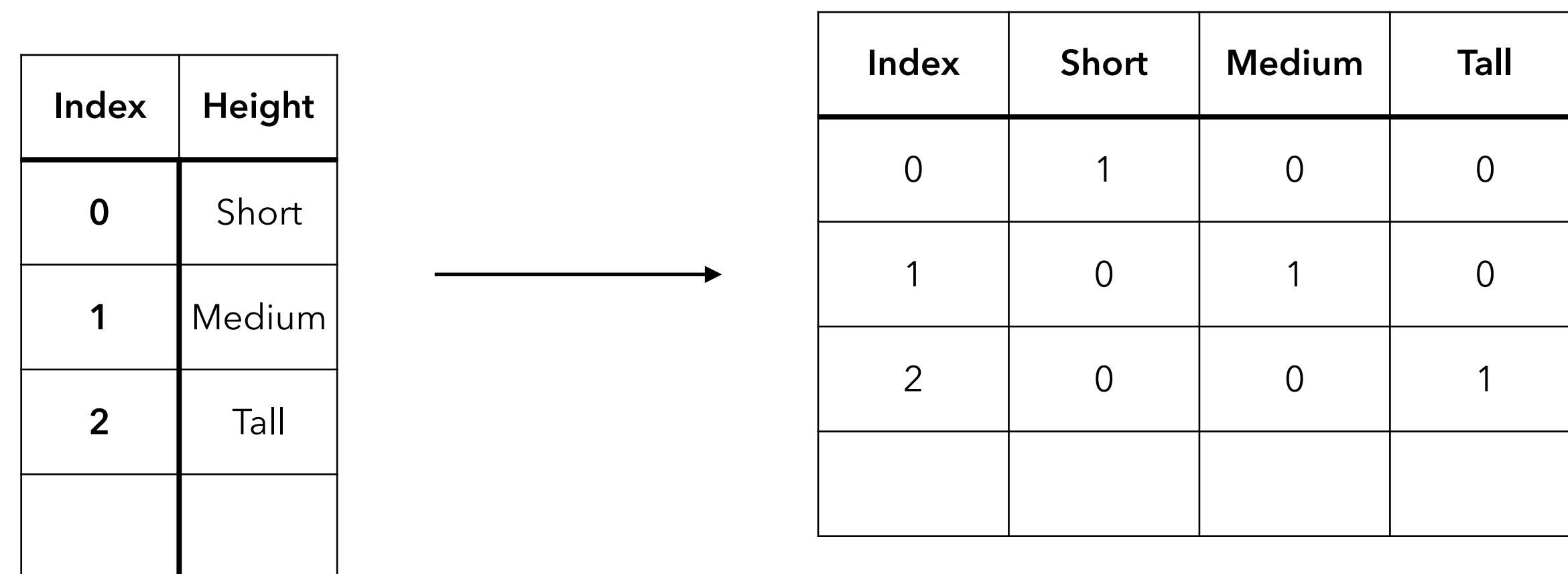


Figure: Label Encoding

Kuhn, M. (2013). Applied Predictive Modeling. Springer

One-Hot Encoding

One-Hot Encoding is a method used to convert categorical variables into a format (Binary Matrix) that can be provided to machine learning algorithms to improve model performance.



The diagram illustrates the process of One-Hot Encoding. On the left, there is a table with two columns: 'Index' and 'Height'. The 'Index' column has values 0, 1, and 2. The 'Height' column has values 'Short', 'Medium', and 'Tall'. An arrow points from this table to a second table on the right. The second table has four columns: 'Index', 'Short', 'Medium', and 'Tall'. The 'Index' column has values 0, 1, and 2. The 'Short' column has values 1, 0, and 0 respectively. The 'Medium' column has values 0, 1, and 0 respectively. The 'Tall' column has values 0, 0, and 1 respectively. This represents the one-hot encoding of the 'Height' variable.

Index	Height
0	Short
1	Medium
2	Tall

Index	Short	Medium	Tall
0	1	0	0
1	0	1	0
2	0	0	1

Figure: One-Hot Encoding

Bishop, C. M. (2006). Pattern recognition and machine learning. Springer google schola, 2, 1122-1128.

Table of Content

- Pre-Alignment
- Theory: Pre-Processing - What is That?
- Case Study: Exploratory Data Analysis

Case Study



STEAM®

What is Steam

- [...] one of the most popular digital game delivery platform (Lin et. al., 2019, p. 1).
- Steam, one of the largest gaming networks for computer games in the world (Baumann et. Al., 2018, p. 1).

Baumann, F., Emmert, D., Baumgartl, H., & Buettner, R. (2018). Hardcore Gamer Profiling: Results from an unsupervised learning approach to playing behavior on the Steam platform. *Procedia Computer Science*, 126, 1289-1297.
 Lin, D., Bezemer, C. P., Zou, Y., & Hassan, A. E. (2019). An empirical study of game reviews on the Steam platform. *Empirical Software Engineering*, 24, 170-207.

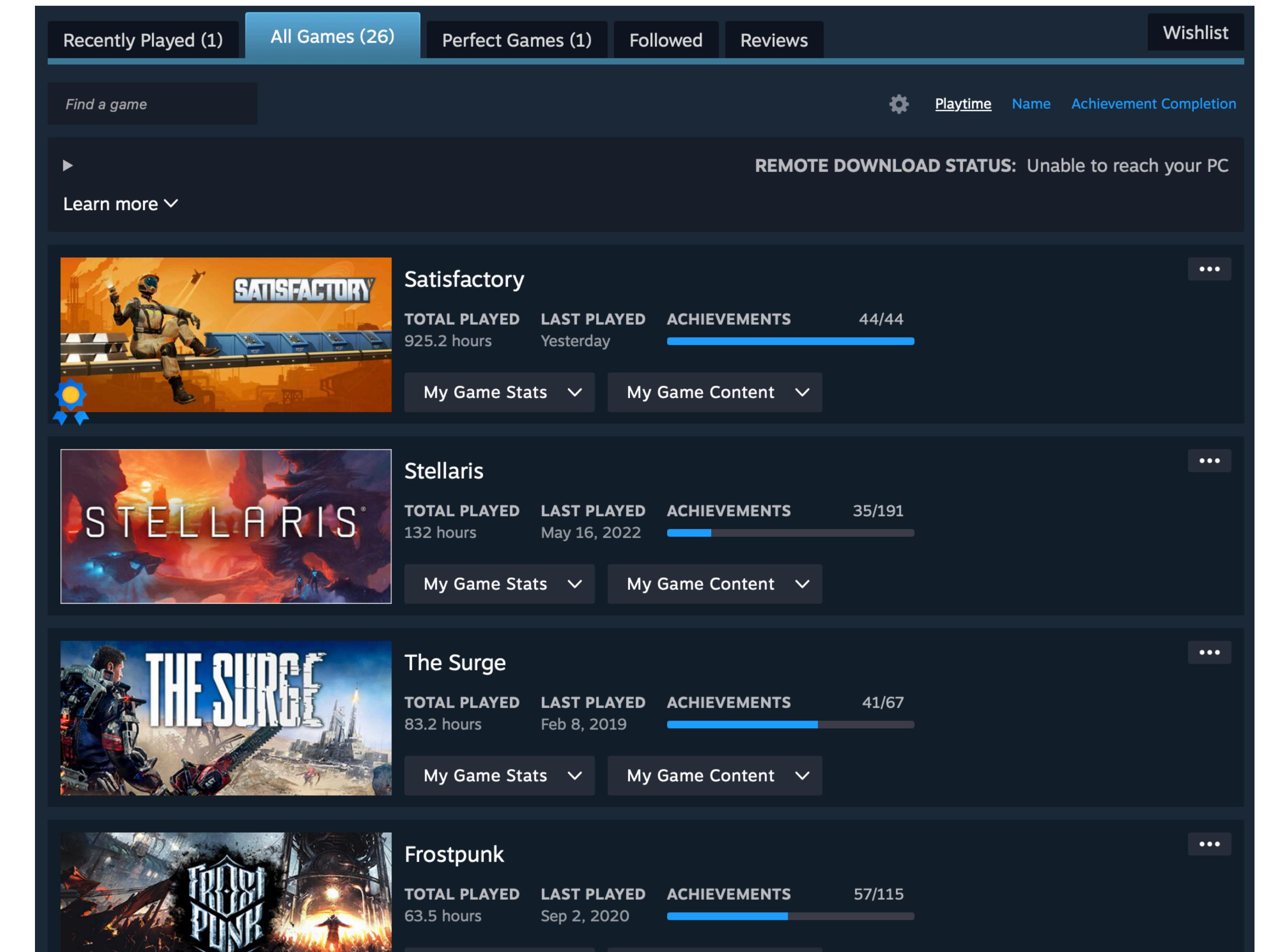
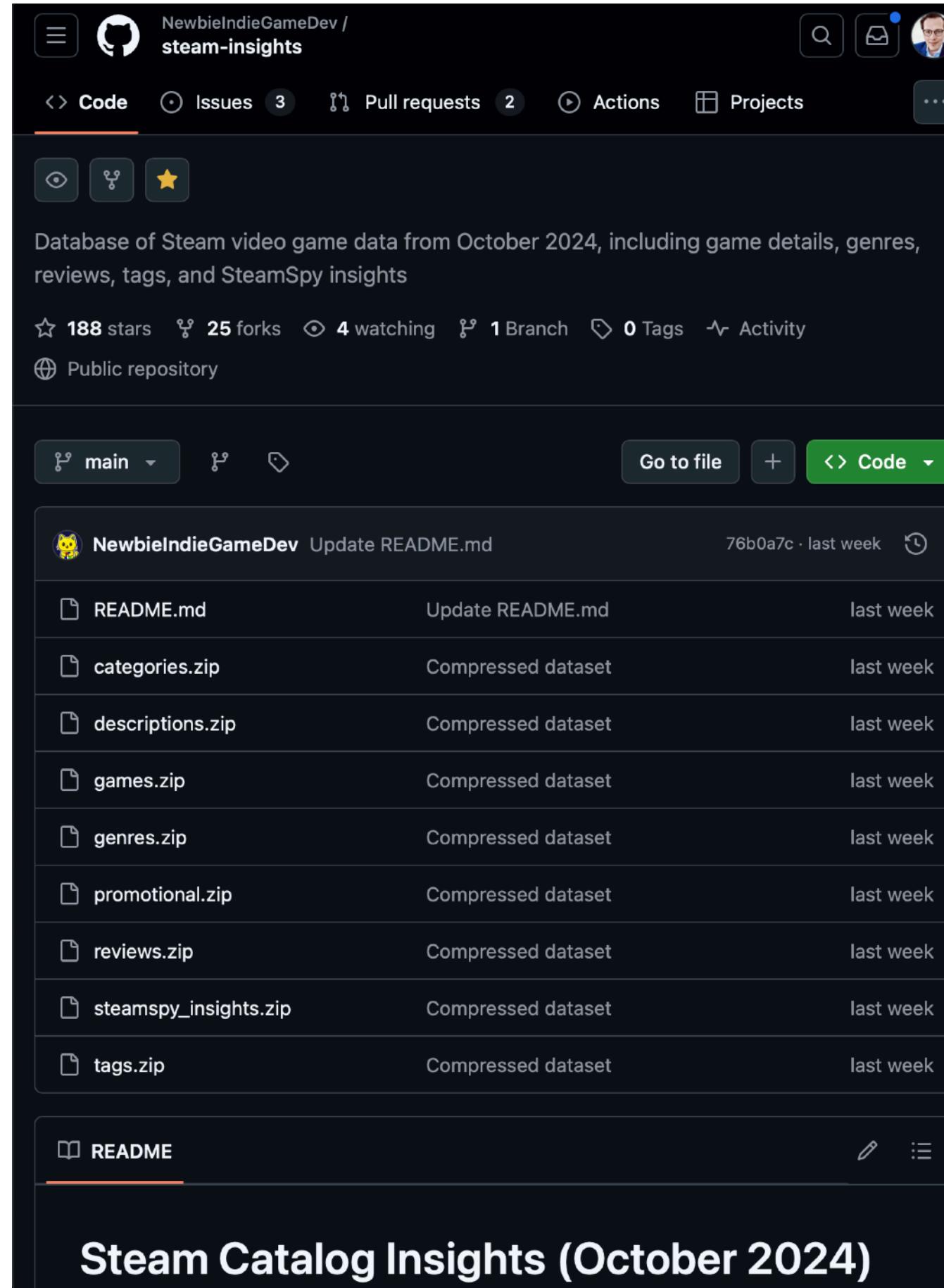


Figure: Steam Licensed Games

Data Sets: High Level View



<https://github.com/NewbieIndieGameDev/steam-insights>

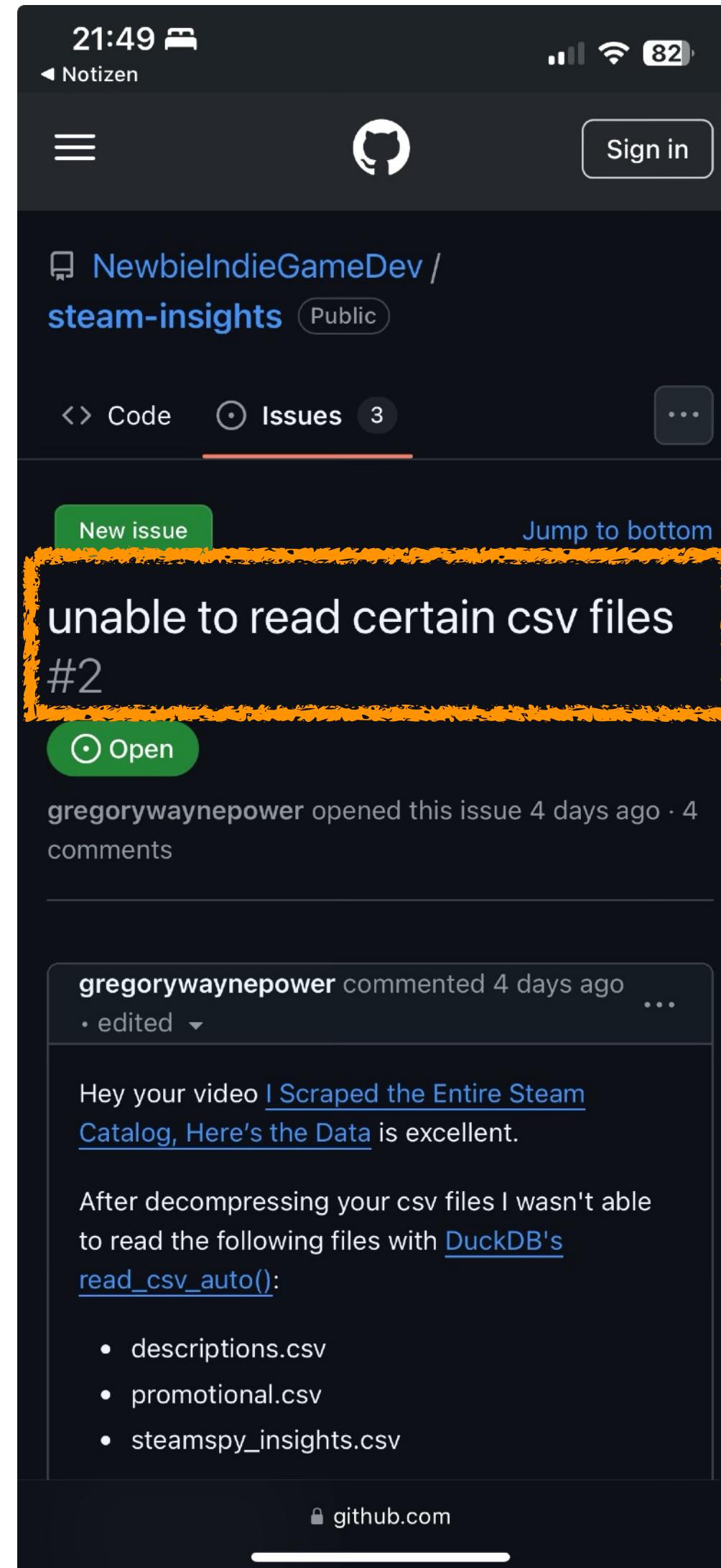
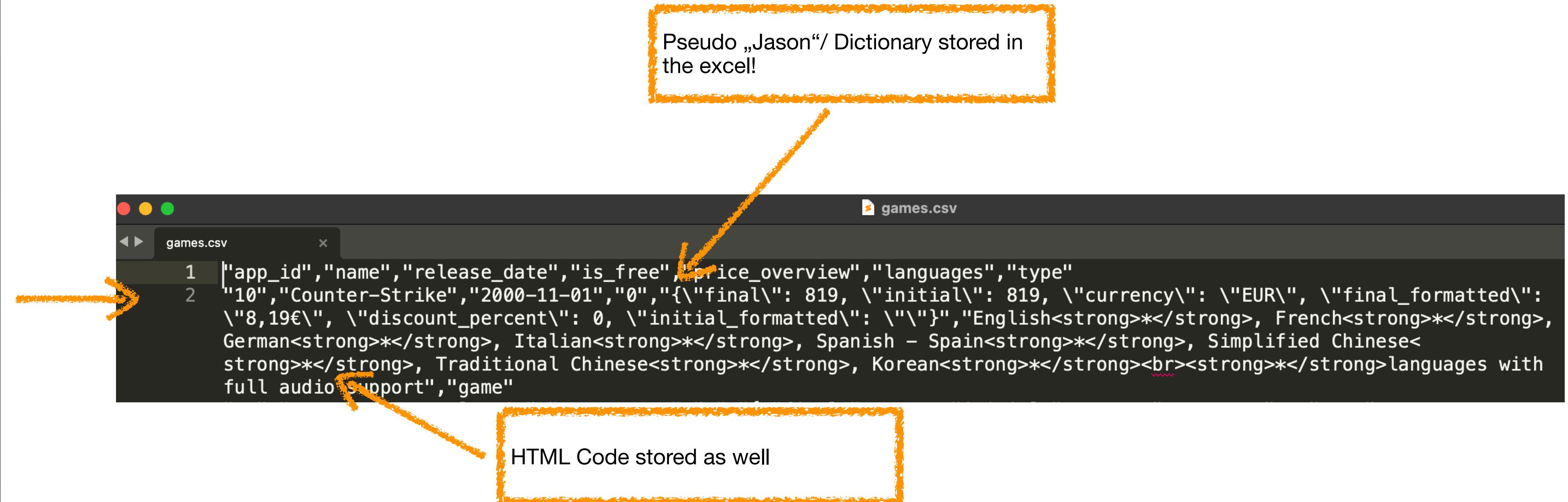
Figure: Git Repo for the Raw Data



<https://www.youtube.com/watch?v=qjNv3qv-YbU>

Figure: Corresponding YouTube Video

Limitation: Data Set is not Perfect

The terminal window displays a CSV file named 'games.csv' with two rows of data. The first row contains column headers: 'app_id', 'name', 'release_date', 'is_free', 'price_overview', 'languages', and 'type'. The second row contains values for these columns. An annotation with a yellow arrow points from the GitHub issue to the terminal window, highlighting the problematic data structure. Another annotation with a yellow arrow points from the terminal window to the right, with the text 'HTML Code stored as well'.

```

1 ["app_id", "name", "release_date", "is_free", "price_overview", "languages", "type"]
2 "10", "Counter-Strike", "2000-11-01", "0", "{\"final\": 819, \"initial\": 819, \"currency\": \"EUR\", \"final_formatted\": \"8,19€\", \"discount_percent\": 0, \"initial_formatted\": \"\"}", "English<strong>*</strong>, French<strong>*</strong>, German<strong>*</strong>, Italian<strong>*</strong>, Spanish – Spain<strong>*</strong>, Simplified Chinese<strong>*</strong>, Traditional Chinese<strong>*</strong>, Korean<strong>*</strong><br><strong>*</strong>languages with full audio support", "game"

```

Do you have any idea why we see this kind of raw data in the .csv?

Figure: CSV Files not Perfect, see Project Issues.

Your Chance

- A) What are the top ten and lowest ten games according to the owners? (Tip: Use owners_conservative.)
- B) Which publishers are the most successful in the dataset based on the number of published games?
- C) How does the number of published games change over time?
- D) What do the prices look like?



Appendix

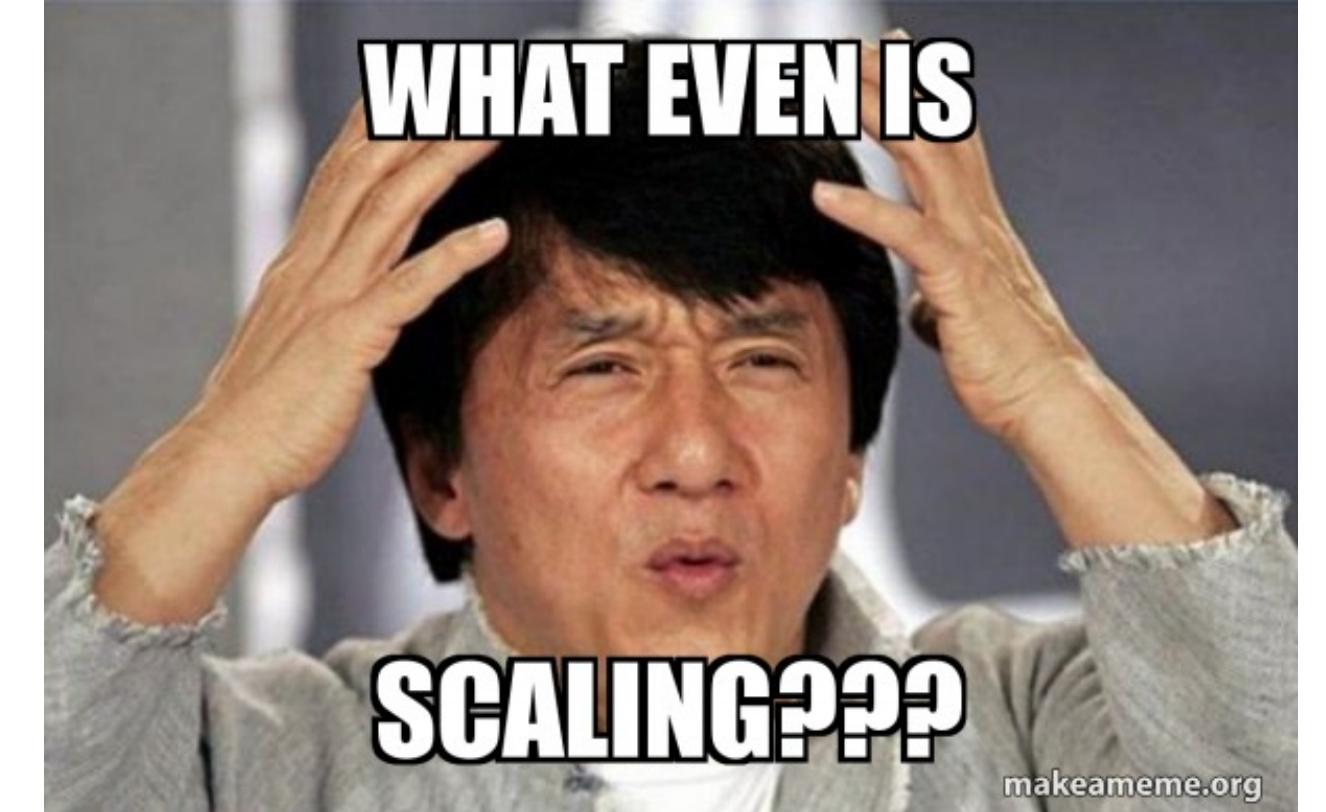
Feature Scaling: Definition

Feature scaling is the process of adjusting the features of a dataset to a uniform scale, ensuring that the differences between values are maintained without distortion.

Normalization: Scale the values of numerical feature into a common scale without distorting differences in the ranges of values or losing information. **Min-Max Scaler:** transforms the features to a specific range, **typically between 0 and 1**.

Standardization, also known as **z-score normalization**, is a scaling method **that centers the values of a feature around its mean, while ensuring a unit standard deviation**. Ready to use is: StandardScaler in sklearn.

Standardization does not always improve the model performance. Effectiveness depends on main working principle of the algorithms. For instance could help by distance calculation, such as K-means, Support Vector Machine (SVM), and Artificial Neural Networks (ANN) to reduce the number space.



<https://makeameme.org/meme/what-even-is-5cda27>

```
from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
X_train_scaled = scaler.transform(X_train)  
X_test_scaled = scaler.transform(X_test)
```

Grandvalet, Y., & Canu, S. (2002). Adaptive scaling for feature selection in SVMs. Advances in neural information processing systems, 15.
Murphy, K. P. (2012). Machine Learning—A probabilistic Perspective. The MIT Press.

Feature Scaling: General Example

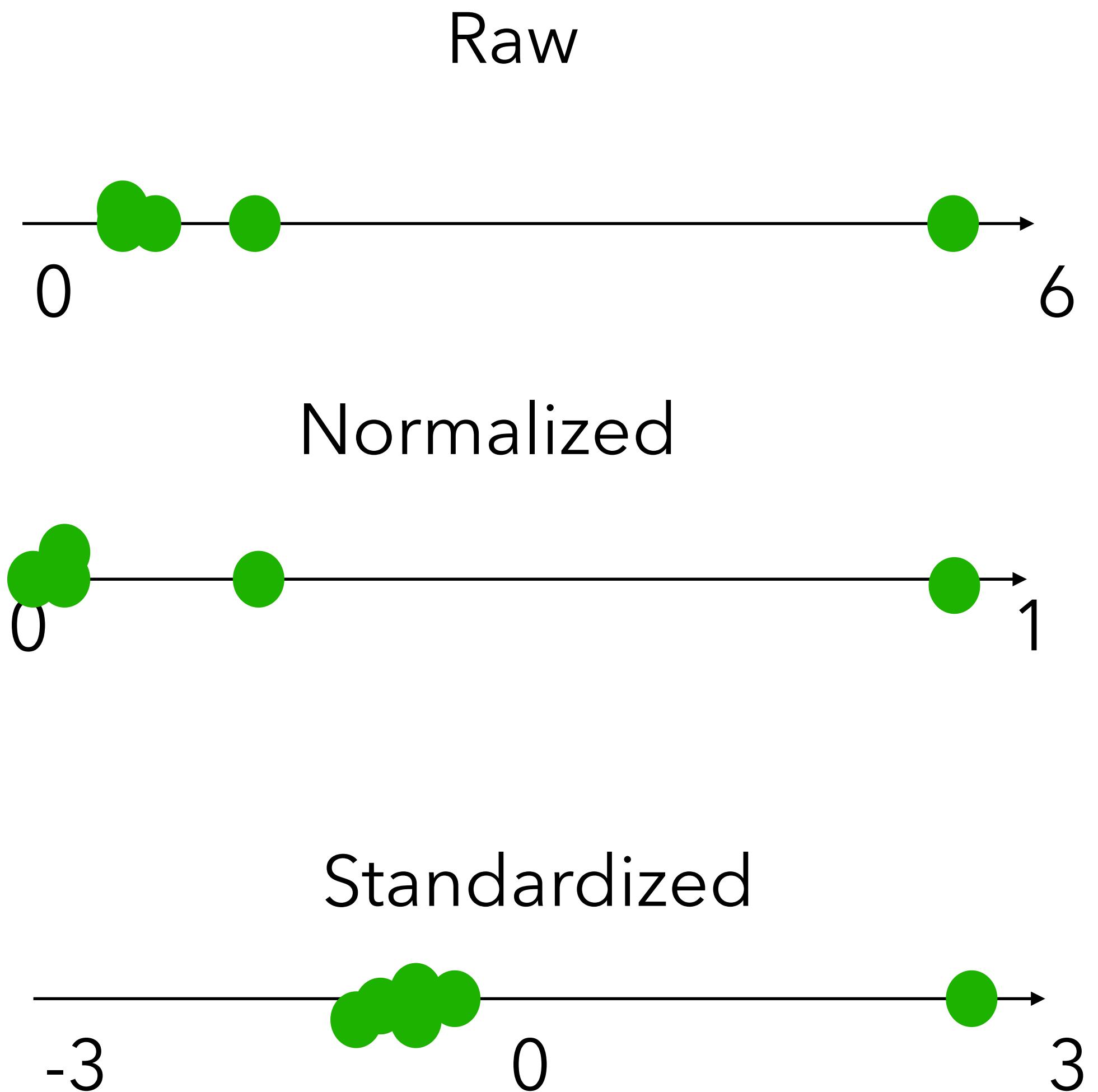


Figure: Illustration of Raw Data, Normalized Data, and Standardized Data

Scaling Method vs. Model: Some Ideas

Scaling Method	Description	When to Use	Research Papers
Normalization (Min-Max Scaling)	Scales features to a fixed range, typically [0, 1].	<ul style="list-style-type: none"> - Neural Networks: Helps speed up convergence and stabilize training. - k-Nearest Neighbors (k-NN): Ensures all features contribute equally to distance calculations. 	<ul style="list-style-type: none"> - LeCun, Y., Bottou, L., Orr, G. B., & Müller, K. (2012). <i>Efficient BackProp</i>. In <i>Neural Networks: Tricks of the Trade</i>. Springer. - Kotsiantis, S. B., Kanellopoulos, D., & Pintelas, P. (2007). <i>Data Preprocessing for Supervised Learning</i>. International Journal of Computer Science and Applications.
Standardization (Z-score Normalization)	Scales features to have a mean of 0 and a standard deviation of 1.	<ul style="list-style-type: none"> - Support Vector Machines (SVM): Improves performance by ensuring features are on the same scale. - Principal Component Analysis (PCA): Helps in accurate principal component computation. - Linear Regression: Improves convergence and model performance. 	<ul style="list-style-type: none"> - Hsu, C.-W., Chang, C.-C., & Lin, C.-J. (2003). <i>A Practical Guide to Support Vector Classification</i>. Technical Report. - Jolliffe, I. T. (2002). <i>Principal Component Analysis</i>. Springer Series in Statistics. - Hastie, T., Tibshirani, R., & Friedman, J. (2009). <i>The Elements of Statistical Learning</i>. Springer.

LeCun, Y., Bottou, L., Orr, G. B., & Müller, K. (2012). Efficient BackProp. In *Neural Networks: Tricks of the Trade* (pp. 9-48). Springer.

Kotsiantis, S. B., Kanellopoulos, D., & Pintelas, P. (2007). Data preprocessing for supervised learning. *International Journal of Computer Science and Applications*, 4(1), 111-117.

Hsu, C.-W., Chang, C.-C., & Lin, C.-J. (2003). A practical guide to support vector classification. Technical Report. Retrieved from <https://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>

Jolliffe, I. T. (2002). *Principal Component Analysis*. Springer Series in Statistics.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.

Outlier Detection and Removal

Outliers are data points that significantly deviate from the rest of the data, and they can arise due to various reasons such as measurement errors, variability in the data, or genuine anomalies.

Note: There is no clear definition of outliers, it depends on the interpretation!

Characteristics of Outliers:

Unusual: Outliers are distinctly different from other data points.

Influential: They can heavily influence statistical analyses, such as mean and variance.

Context-Dependent: What constitutes an outlier can vary depending on the context of the data and the specific analysis being performed.

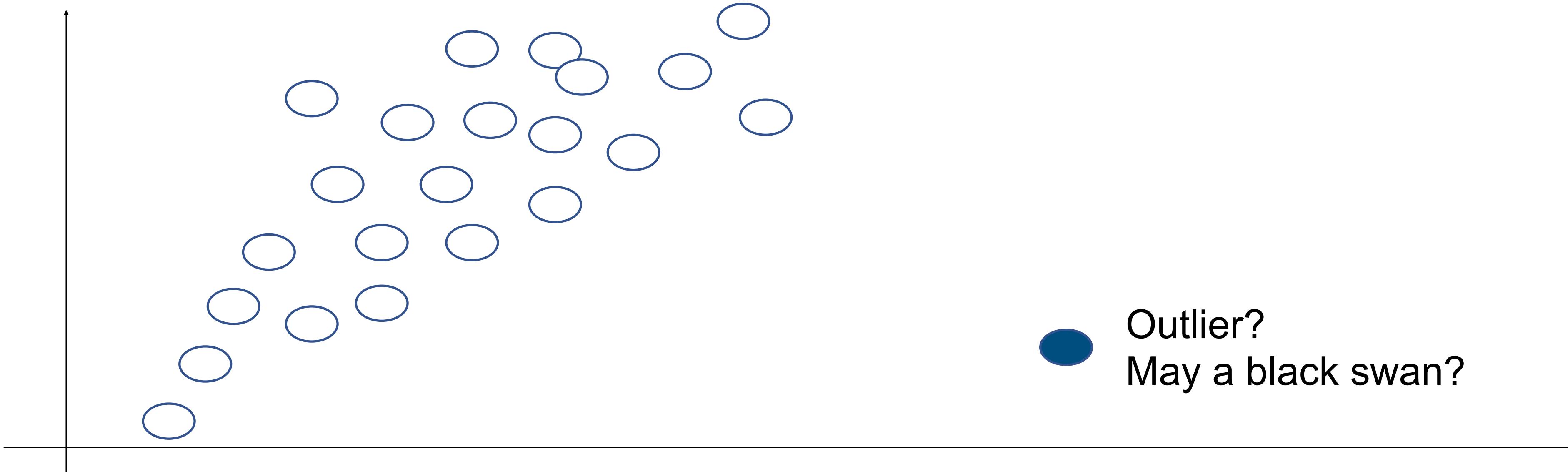


Figure: The Problems with Outlier Detection

Hodge, V. J., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2), 85-126.

Find Outliers: Have a Look

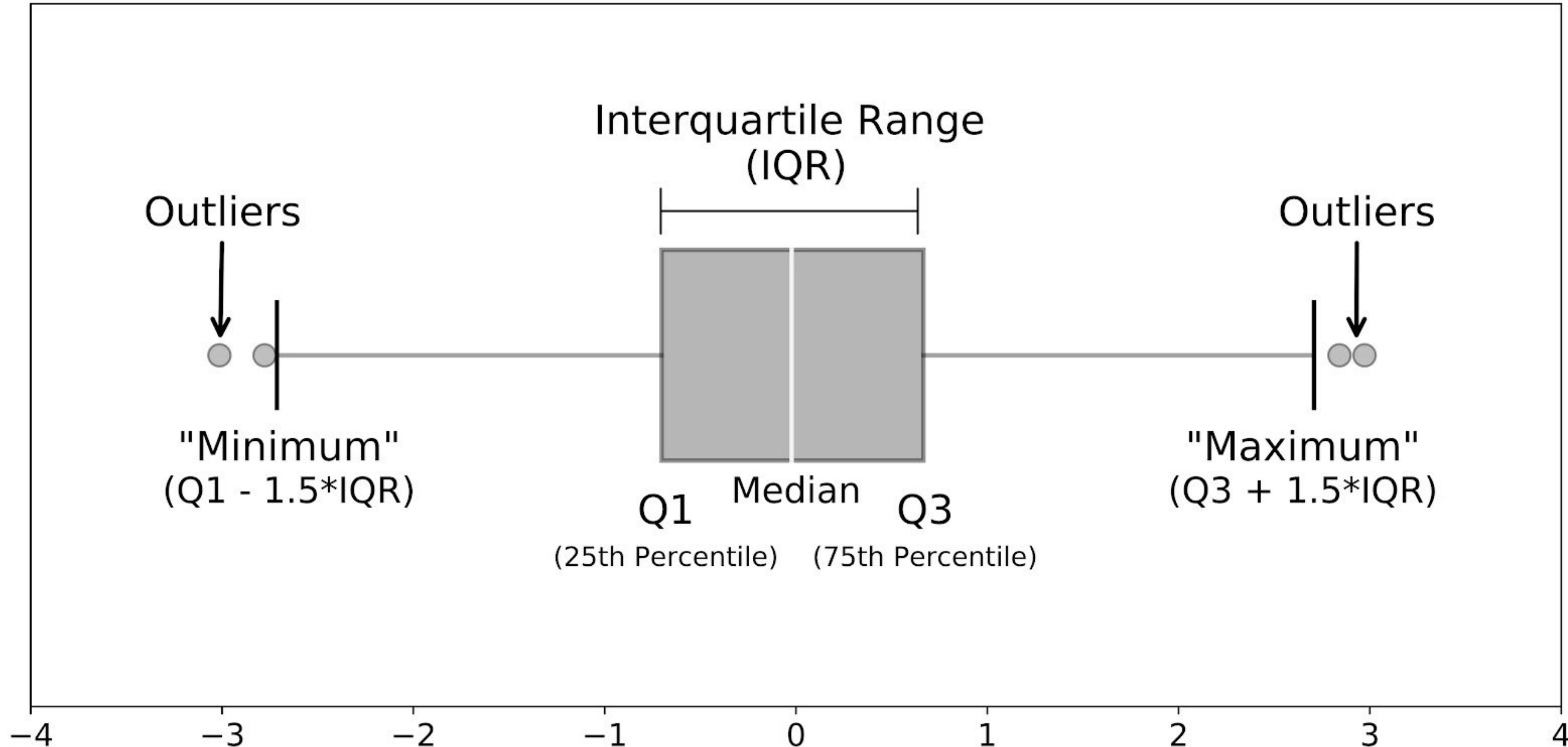
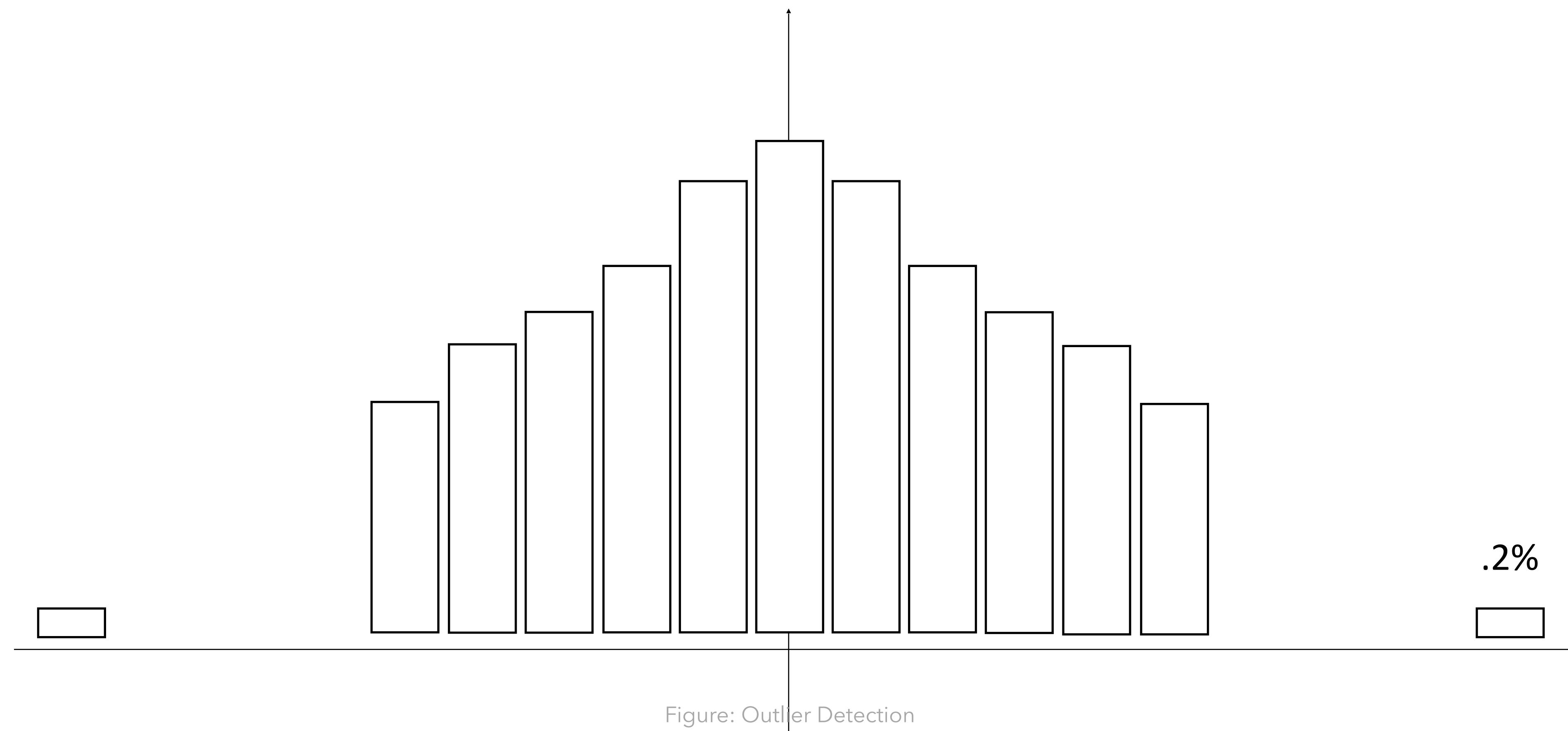


Figure: Boxplot Example

Example Approach

Calculation: mean value +/- 1.98 of the standard deviation



Outlier Detection and Removal

- Statistical Methods:
 - **Z-Score:** Measures how many standard deviations a data point is from the mean of the data. A high absolute Z-score indicates an outlier.
 - **Interquartile Range (IQR):**
 - Defines outliers as those lying below $Q1 - 1,5 * IQR$ or $Q3 + 1,5 * IQR$ above
- Machine Learning-Based Methods: K-NN, SVM ...
- Visual Methods:
 - Boxplot
 - Scatterplot

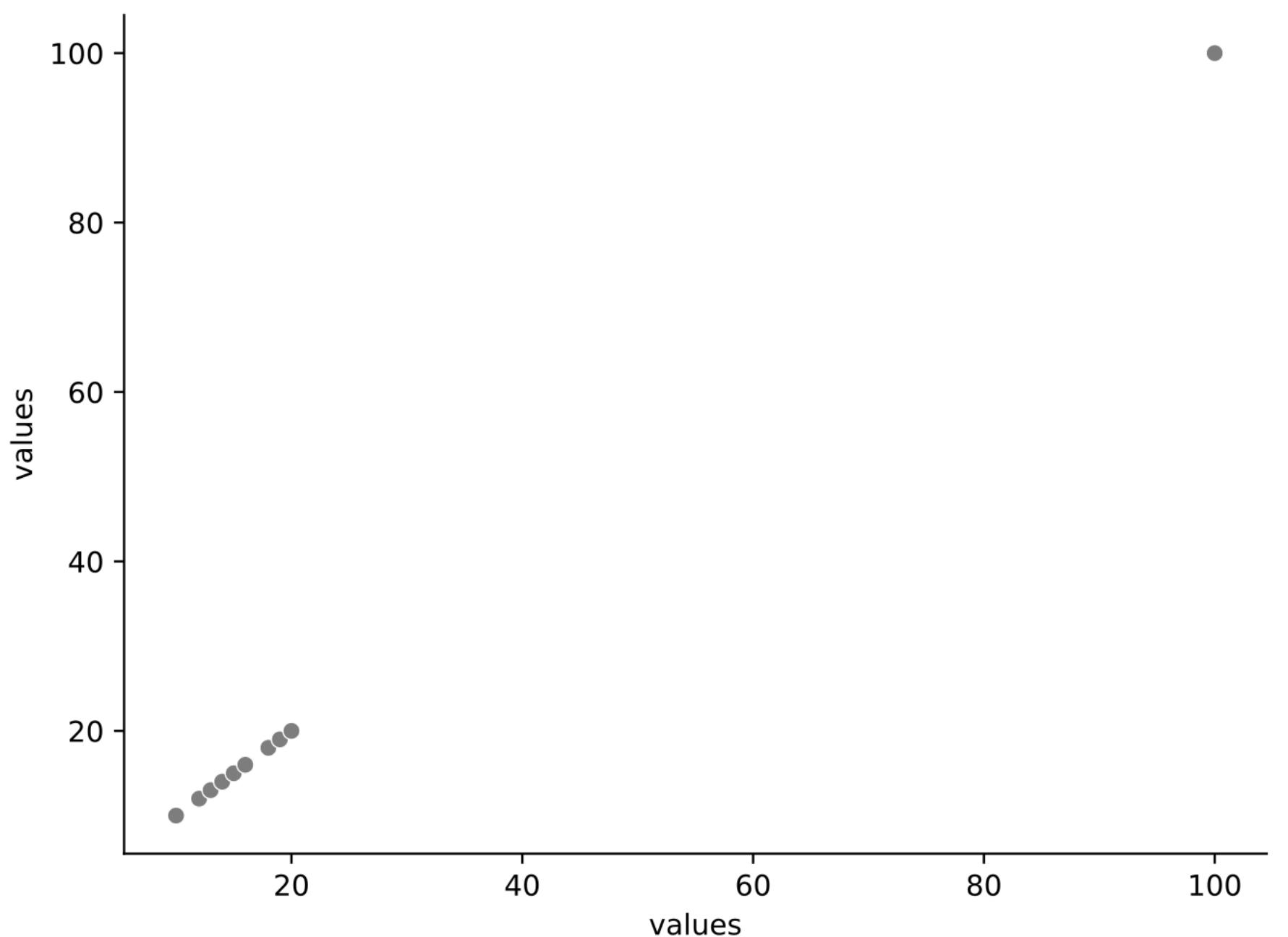


Figure: Outlier Detection

Outlier Detection: Limitation

- **Context Dependency:** The effectiveness of outlier detection methods often depends on the context of the data and the specific application, making it necessary to understand the data characteristics before choosing a method.
- **Statistical Methods:** May assume normality or symmetry and can be sensitive to outliers in the computation of summary statistics.
- **Machine Learning-Based Methods:** May require careful parameter tuning and can be computationally expensive.
- **Distance-Based Methods:** Sensitive to parameter choices and may struggle with high-dimensional data.
- **Visual Methods:** Limited to specific scenarios and can be subjective in interpretation.

Alghushairy, O., Alsini, R., Soule, T., & Ma, X. (2020). A review of local outlier factor algorithms for outlier detection in big data streams. *Big Data and Cognitive Computing*, 5(1), 1.

Barbato, G., Barini, E. M., Genta, G., & Levi, R. (2011). Features and performance of some outlier detection methods. *Journal of Applied Statistics*, 38(10), 2133-2149.

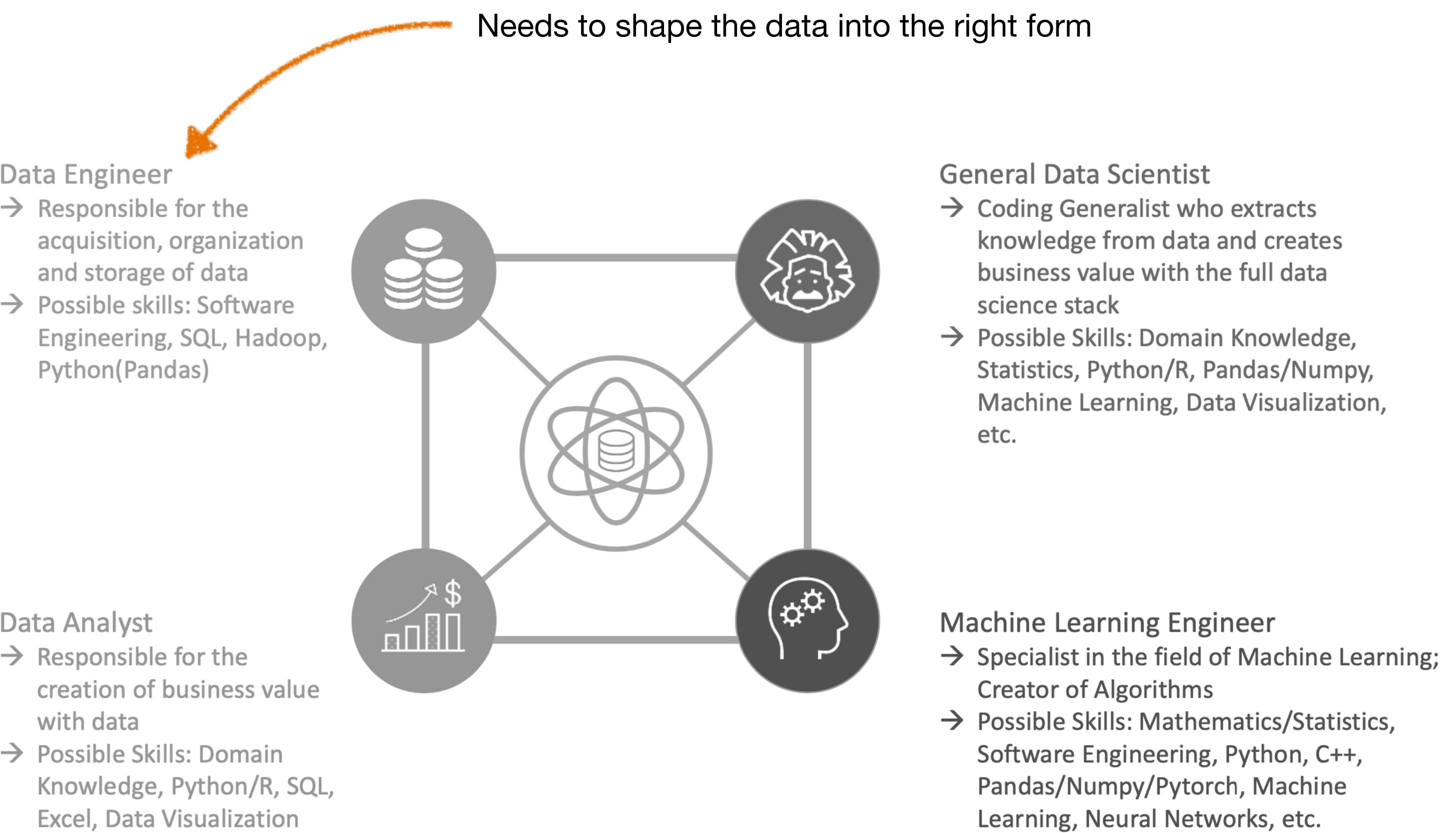
Hodge, V., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial intelligence review*, 22, 85-126.

From Data to Machine Learning ... Easy? No!

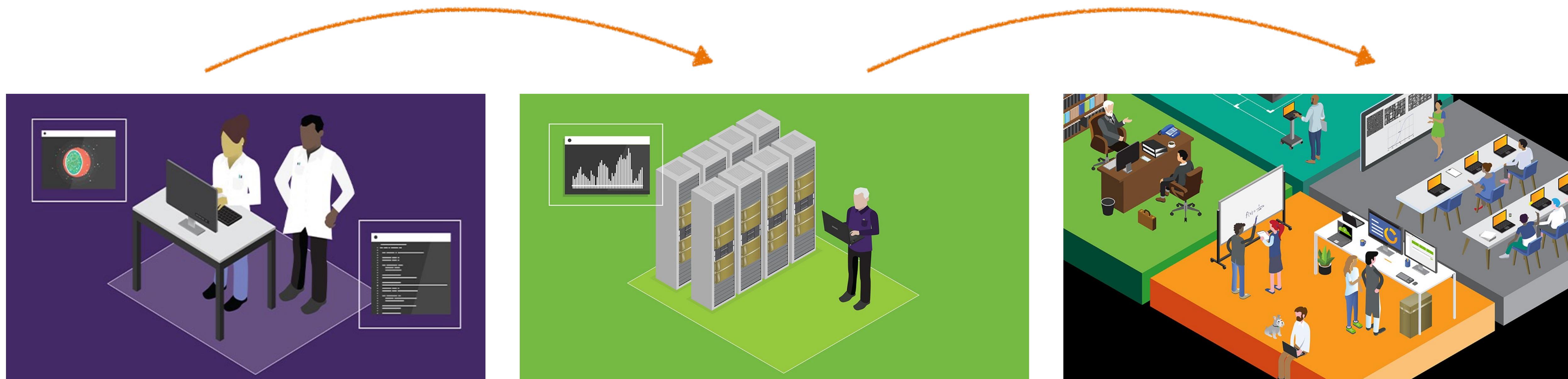
Database Management	Machine Learning
Database is an active, evolving entity	Database is just a static collection of data
Records may contain missing or erroneous information	Instances are usually complete and noise-free
A typical field is numeric	A typical feature is binary
A database typically contains millions of records	Data sets typically contain several hundred instances
AI should get down to reality	"Databases" is a solved problem and is therefore uninteresting

Table 2. Conflicting Viewpoints between [1] Database Management and Machine Learning.

[1] Frawley, W. J., Piatetsky-Shapiro, G., & Matheus, C. J. (1992). Knowledge discovery in databases: An overview. *AI magazine*, 13(3), 57-57.



Pre-Processing: Challenge Workflow



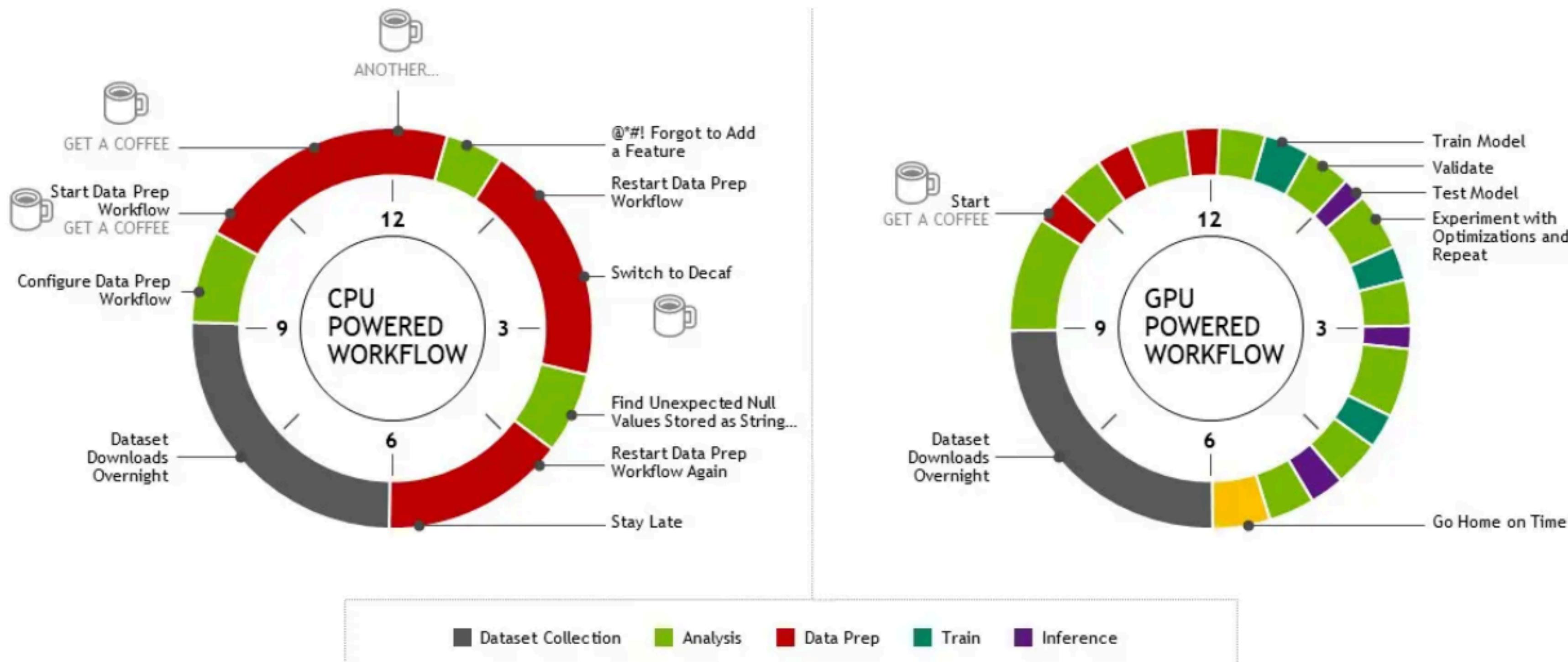
„Local“ Development

Server Deployment

Use in Production

Pre-Processing: Critical Reflection

DAY IN THE LIFE OF A DATA SCIENTIST



Therefore, it is important to realize that it is common practice to dedicate about 80% of the labor and time within the KDD process to data preparation, as a representative survey on a popular data science portal shows. <https://bit.ly/2WwVPho>.

In: Abdel-Karim, B. M., Pfeuffer, N., & Hinz, O. (2021). Machine learning in information systems-a bibliographic review and open research issues. *Electronic Markets*, 31(3), 643-670

<https://blogs.nvidia.com/blog/2018/11/15/accelerated-data-science-hpc/>

Domain Knowledge: Why is this Important?

Domain knowledge can be define as knowledge about the environment in which the target implementation operates [1].



Data Science Team

[1] Hjørland, B. & Albrechtsen, H. (1995). [Toward A New Horizon in Information Science: Domain Analysis](#). Journal of the American Society for Information Science, 1995, 46(6), p. 400–425.



Images:
<https://www.theforage.com/blog/careers/investment-banking>
<https://unsplash.com/de/s/fotos/Geschäft>
<https://stock.adobe.com/de/images/nerd-at-the-computer-side-view-of-young-nerd-man-working-at-the/55780370>