



Data Analysis Fundamentals

Lecture 3: Understanding Data Operations

Vikas Kumar
Indian Institute of Technology Bombay

07th September 2023

Outline

- Recap
- Data Operations
- Select, Filter - Simple vs Complex, Sort, Group & Aggregate
- Merge, Pivot, Unpivot, Windowing

Recap

- Data and its evolution
- Data Science as a career option and typical job roles
- Data Science Pipeline and toolkit of a Data Analyst
- Excel and its history
- Data, Binning and Granularity
- Metrics and KPIs

Data Operations

- Processes applied to the data
- Processes can be mathematical/non-mathematical



`=SUM(A1:A5)`
`=COUNT(A1:A5)`
`=AVERAGE(A1:A5)`
`=MIN(A1:A5)`
`=MAX(A1:A5)`
`=IF(A1>33,"P","F")`

Dimensions and Facts

- A dimension is a measure of a physical variable (without numerical values).
- A unit is a way to assign a number or measurement to that dimension.
- For example, length is a dimension, but it is measured in units of feet (ft) or meters (m).
- A numeric attribute for which data operations can be done is called fact.
- A fact is the numerical part of the dimension.

Cardinality

- Cardinality means how the entities are arranged to each other.
- The relationship structure between entities in a relationship set.

one-to-one (1:1)



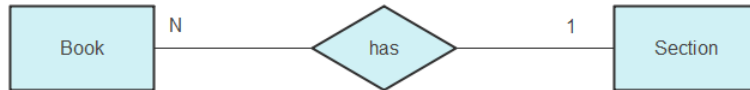
✓ One student can have only one student ID

one-to-many (1:N)



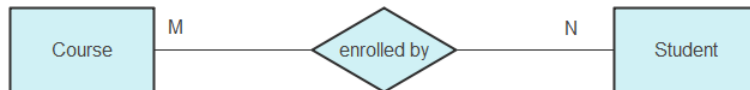
✓ Edyoda has multiple trainers

many-to-one (N:1)



✓ One trainer teaches multiple students in the class

many-to-many (M:N)



✓ One trainer takes several modules, and same module is taken by different trainers

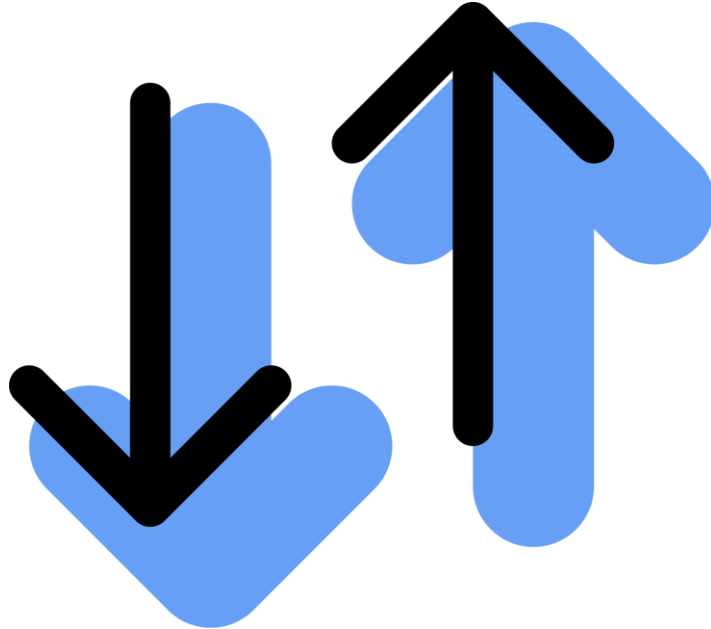
Filter

- Data filtering is the process of examining a dataset to exclude, rearrange, or apportion data according to certain criteria.
- Process of choosing a smaller part of your dataset and using that subset for viewing or analysis.



Sort

- Data sorting is any process that involves arranging the data into some meaningful order to make it easier to understand, analyze or visualize.
- Sorting refers to ordering data in an increasing or decreasing manner.



Group

- Grouped data are data formed by aggregating individual observations of a variable into groups, so that a frequency distribution of these groups serves as a convenient means of summarizing or analyzing the data.

GROUPED DATA VS UNGROUPED DATA

- ◉ **Ungrouped data** – Data that has not been organized into groups. Also called as raw data.

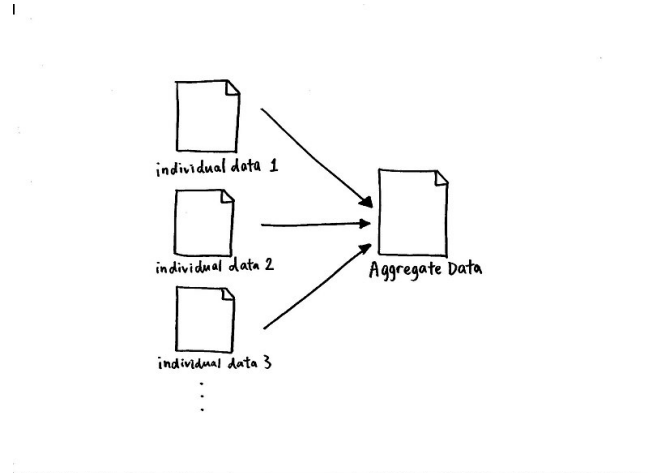
Data	Frequency
2	8
3	4
5	6
7	7
8	2
9	5

- ◉ **Grouped data** - Data that has been organized into groups (into a frequency distribution).

Data	Frequency
2 - 4	5
5 - 7	6
8 - 10	10
11 - 13	8
14 - 16	4
17 - 19	3

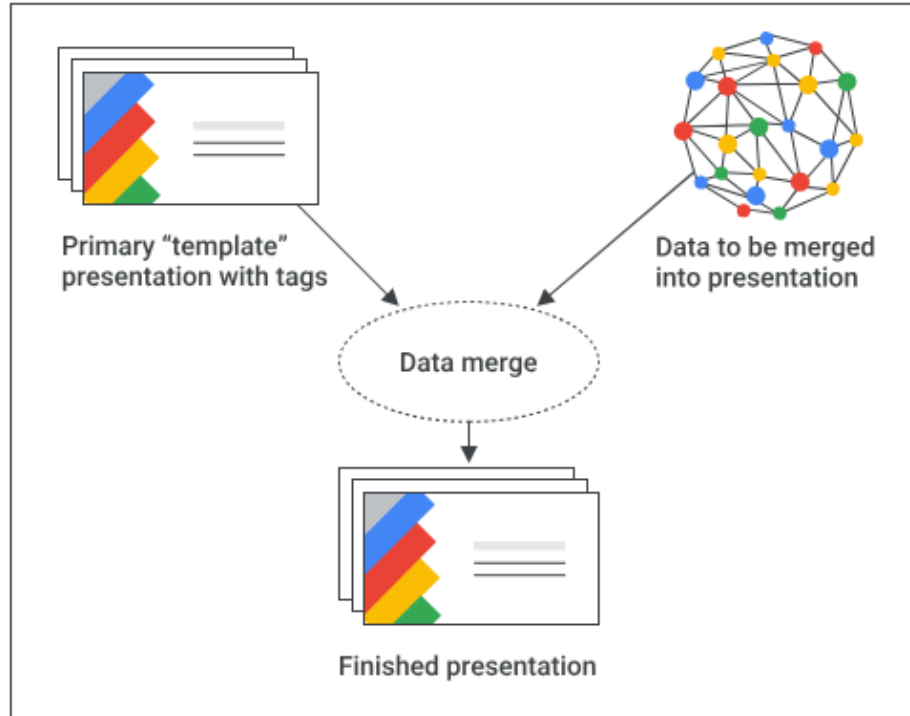
Aggregate

- Data aggregation is the process where raw data is gathered and expressed in a summary form for statistical analysis.
- Aggregate data is high-level data which is acquired by combining individual-level data.
- AVERAGE, SUM, PRODUCT, COUNT, COUNTA, MAX, or MIN



Merge

- Data merging is the process of combining two or more similar records into a single one.



Windowing

- Window functions are sometimes used in the field of statistical analysis to restrict the set of data being analyzed to a range near a given point.

A word cloud featuring the phrase "Thank You" in numerous languages and colors. The central and largest text is "thank you" in red. Other prominent words include "gracias" in green, "danke" in blue, "merci" in orange, and "shukriya" in purple. Smaller words in various colors include "arigatō", "dank je", "teşekkür ederim", "ngiyabonga", "shukra jayila", "spas", "barka", "welalin", "tack", "misaotra", "matondo", "paldies", "grazzi", "mahalo", "tapadh leat", "xhala", "asante", "manana", "tenki", "murekaze", "chokran", "mamnun", "dyaquyo", "mochchakkeram", "go raibh maith agat", "arigatō", "dakujem", "trugarez", "merci", "shukriya", "merce", "merci", "diolch", "dhanyavadagalu", "tanemirt", "rahmet", "xixie", "eucharistw", "sagolun", "sukriya", "kop khun krap", "sulpay", "gracias ago", "gracies", "chnorakaloutioun", "najis tuke", "kam sah hamnida", "rahat", "tomasake", "banyabad", "didi", "madloba", "mes", "dekuji", "sobodi", "obrigado", "dziękuje", "hvala", "maururu", "kösönöm", "bayatalaa", "gracie", "dhanyavad", "kiitos", "dankie", "faafetai lava", "spasibo", "Баярлалаа", "рахмат", "danke", "謝謝", "ngiyabonga", "shukra jayila", "spas", "barka", "welalin", "tack", "misaotra", "matondo", "paldies", "grazzi", "mahalo", "tapadh leat", "xhala", "asante", "manana", "tenki", "murekaze", "chokran", "mamnun", "dyaquyo", "mochchakkeram", "go raibh maith agat", "arigatō", "dakujem", "trugarez", "merci", "shukriya", "merce", "merci", "diolch", "dhanyavadagalu", "tanemirt", "rahmet", "xixie", "eucharistw", "sagolun", "sukriya", "kop khun krap", "sulpay", "gracias ago", "gracies", "chnorakaloutioun", "najis tuke", "kam sah hamnida", "rahat", "tomasake", "banyabad", "didi", "madloba", "mes", "dekuji", "sobodi", "obrigado", "dziękuje", "hvala", "maururu", "kösönöm", "bayatalaa", "gracie", "dhanyavad", "kiitos", "dankie", "faafetai lava", "spasibo", "Баярлалаа", "рахмат", "danke", "謝謝".