

Data Cleaning

Freitag, 11. Juli 2025 10:29

Dirty Data

- ↳ incomplete, incorrect or irrelevant data
- ↳ to the problem we're trying to solve

Clean Data

- ↳ is complete, correct and relevant
- ↳ to the problem we're trying to solve

Data Engineers

- ↳ transform data into a useful format for analysis
- ↳ and give it a reliable infrastructure

Data Warehousing Specialists

- ↳ develop processes and procedures
- ↳ to effectively store and organise data

Null

- ↳ indication that a value does not exist in a dataset
- ↳ it means that the customer skipped that question
- ↳ it is not a 0 - a zero can be a zero as a response

Field

- ↳ a single piece of information from a row or a column of a spreadsheet

Field Length

- ↳ a tool for determining how many characters can be keyed into a field

Data Validation

- ↳ a tool for checking the accuracy and quality of data
- ↳ before adding or importing it
- ↳ a form of data cleansing

Types of dirty data



Duplicate data



Outdated data



Incomplete data



Incorrect/inaccurate data



Inconsistent data

Duplicate Data

Description	Possible causes	Potential harm to businesses
Any data record that shows up more than once	Manual data entry, batch data imports, or data migration	Skewed metrics or analyses, inflated or inaccurate counts or predictions, or confusion during data retrieval

Outdated Data

Description	Possible causes	Potential harm to businesses
Any data that is old which should be replaced with newer and more accurate information	People changing roles or companies, or software and systems becoming obsolete	Inaccurate insights, decision-making, and analytics

Incomplete Data

Description	Possible causes	Potential harm to businesses
Any data that is missing important fields	Improper data collection or incorrect data entry	Decreased productivity, inaccurate insights, or inability to complete essential services

Datei auswählenKeine auswählen

Incorrect/Inaccurate Data

Description	Possible causes	Potential harm to businesses
Any data that is complete but inaccurate	Human error inserted during data input, fake information, or mock data	Inaccurate insights or decision-making based on bad information resulting in revenue loss

Inconsistent Data

Description	Possible causes	Potential harm to businesses
Any data that uses different formats to represent the same thing	Data stored incorrectly or errors inserted during data transfer	Contradictory data points leading to confusion or inability to classify or segment customers

Dirty Data has huge impact on companies across various fields!

Data Integrity [Rep]

Freitag, 11. Juli 2025 10:29

Data Integrity

- ↳ ensuring that data is correct
- ↳ Repetition - 4 parts

Validity

- ↳ measures conform to defined business rules or constraints

Validity

Definition

The concept of using data integrity principles to ensure measures conform to defined business rules or constraints

Example

Data collected five years ago used technology that is not approved or supported by the business

Accuracy

- ↳ compared to reliable outside sources

Accuracy

Definition

The degree of conformity of a measure to a standard or a true value

Example

Addresses in the business database are identified as incorrect when compared to the public postal service database

Completeness

- ↳ ensuring that all important data is included

Completeness

Definition

The degree to which all required measures are known

Example

NULL/missing value for the item "Number of employees per store"

Consistency

- ↳ measures are formatted or structured
- ↳ the same way across systems

Consistency

Definition

The degree to which a set of measures is equivalent across systems

Example

Date of store opening stored in both MM/DD/YYYY and MM/YY formats

Datei auswählen
Keine ausgewählt

Data Merging / Errors

Freitag, 11. Juli 2025 10:29

Merger

- ↳ an agreement that unites two organisations into a single new one
- ↳ in case you have data from two sources

Data Merging

- ↳ process of combining two or more datasets into a single dataset

Compatibility

- ↳ how well two or more datasets are able to work together

How To Access Compatibility

- ↳ Do I have all the data I need?
- ↳ Does the data I need exist within these datasets?
- ↳ Do the data need to be cleaned, or are they ready for me to use?
- ↳ Are the datasets cleaned to the same standard?

Errors While Cleaning Data



Common mistakes to avoid

- ↳ Not checking for spelling errors
- ↳ Forgetting to document errors (*can be a big time safer*)
- ↳ Not checking for misfielded values
- ↳ Overlooking missing values
- ↳ Only looking at a subset of the data (*every dataset equal attention*)
- ↳ Losing track of business objectives (*new discoveries while cleaning*)
- ↳ Not fixing the source of the error (*fixing same error again and again*)
- ↳ Not analysing the system prior to data cleaning
 - ↳ learn the source of dirty data
- ↳ Not backing up your data prior to data cleaning (*create backups*)
- ↳ Not accounting for data cleaning in your deadlines/process
 - ↳ invest enough time into data cleaning! (*but effectively*)

10 Google Workspace Tips to clean up Data

<https://support.google.com/a/users/answer/9604139?hl=en#zippy=>

Datei auswählen
Keine ausgewählt

Cleaning Spreadsheets 1

Freitag, 11. Juli 2025 10:29

Cleaning Data With Blank Fields

- ↳ apply a filter on every column
- ↳ where it only shows empty cells
- ↳ go column after column
- ↳ delete all empty cells

How To Transpose Data

- ↳ from long format to wide format
- ↳ copy the whole table
- ↳ right-click and "*Paste-Special*"
- ↳ and then choose "*Transpose*"

Get Rid of Extra Spaces in cells

- ↳ select the whole sheet
- ↳ select "*Data*"
- ↳ select "*Data clean-up*"
- ↳ select "*Trim Whitespace*"
- ↳ optional: =TRIM(A1)
- ↳ and then drag auto-fill square

Add-On: Change lower/uppercase/Proper case text

- ↳ process string data
- ↳ Add-On is the easiest way to clean up string data
- ↳ All Uppercase (*makes all letters BIG*)

Delete All Formatting

- ↳ highlight all Rows 1 - 8
- ↳ select the tab "*Format*"
- ↳ select "*Clear Formatting*"

Datei auswählen
Keine ausgewählt

Cleaning Spreadsheets 2

Freitag, 11. Juli 2025 10:29

Conditional Formatting

↳ changes how cells appear when values meet specific conditions

Remove Duplicates

↳ tool that automatically searches for duplicates
↳ and deletes duplicate entries from the spreadsheet

Text String

↳ a group of characters within a cell
↳ mostly composed of letters

Split

↳ tool that divides text around a specific character
↳ puts each fragment into a new, separate cell
↳ e.g. commas recognised as delimiters
↳ delimiter = specified text separator

Concatenate

↳ function joining multiple text strings into a single string

Auto Data Format Setting

↳ mark the column
↳ select "Format"
↳ select "Number"
↳ select "Date"

Splitting Data into Columns

↳ when column has commas or similar
↳ you can automatically split into columns
↳ mark column
↳ select the tab "Data"
↳ select "split into columns"
↳ there is a **Function!**
↳ =SPLIT(Cell, "Split-Sign")
↳ Split-Sign is what determines the place to split

Fixing Errors?! Using "Splitting Data into Columns"

↳ there was a #VALUE! Error
↳ a number in a field had this "717" format
↳ and the splitting function just removed the " "
↳ and solved the error

Datei auswählen
Keine ausgewählt

Cleaning Spreadsheets 3

Freitag, 11. Juli 2025 10:29

Function

- ↳ set of instructions that performs a specific calculation
- ↳ using the data in a spreadsheet

Syntax

- ↳ predetermined structure
- ↳ including all required information
- ↳ and it's proper placement
- ↳ what I always write about function before e.g.

The COUNTIF Function

- ↳ write into an empty cell
- ↳ =COUNTIF (Range:Range; "Value")
- ↳ e.g. =COUNTIF(I2:I72, "<100")
- ↳ count how many cells have value under 100
- ↳ can be used to clean data
- ↳ can be used to check data integrity

The LEN Function

- ↳ if certain piece of information has certain length
- ↳ =LEN(Cell) , e.g. =LEN(A2)
- ↳ tells me how many numbers/letters cell is made of
- ↳ spreading it all over a column can help to find mistakes

The LEFT/RIGHT/MID Function

- ↳ can return a segment of a cell
- ↳ starting from the left, right or from the middle
- ↳ =LEFT(Cell, Number Of Letters)
- ↳ e.g. =LEFT(A2, 5) ; =RIGHT(A2, 4)
- ↳ =MID(Cell, Starting Point Letter Number, Number Of Letters)
- ↳ e.g. =MID(A2, 3, 2)

The CONCATENATE Function

- ↳ joins together two or more text strings
- ↳ put info from two cells into one cell
- ↳ =CONCATENATE(Cell 1, Cell 2)
- ↳ e.g. =CONCATENATE(H2, I2)

The TRIM Function

- ↳ removing leading, trailing and repeated spaces in data
- ↳ if there are too many "spaces" or "-" or ".."
- ↳ =TRIM(Cell), e.g. =TRIM(A2)

Datei auswählen
Keine ausgewählt

Cleaning Spreadsheets 4

Freitag, 11. Juli 2025 10:29

Sorting

- ↳ arranging data into a meaningful order
- ↳ to make it easier to understand, analyse and visualise

Filtering

- ↳ showing only the data that meets a specific criteria
- ↳ while hiding the rest

Pivot Table

- ↳ data summarisation tool used in data processing

VLOOKUP

- ↳ means virtual look up

Data Mapping

- ↳ process of matching fields from one databases to another
- ↳ critical for data migration, integration and more

Schema

- ↳ way of describing how something is organised

Primary Key

- ↳ references a column in which each value is unique

Foreign Key

- ↳ a field within a table that is a primary key in another table

Pivot Tables

- ↳ data summarisation tool
- ↳ used in data processing and data cleaning
- ↳ we used it to copy 2 rows into a table
- ↳ go to the tab "Insert"
- ↳ click on "Pivot Table"
- ↳ there is an self-explaining editor

The VLOOKUP Function

- ↳ vertically searches for a certain value in a column
- ↳ return a corresponding piece of information
- ↳ but usually we're searching across multiple sheets/databases
- ↳ that's what we will do now
- ↳ =VLOOKUP(data to look up, 'where to look'!Range, column, false)
- ↳ =VLOOKUP(A2, 'Sheet 2'!A1:B31, 2, false)
- ↳ **take value in A2** of Sheet 1
- ↳ **checking value in Sheet 2** from A1 to B31
- ↳ we **want the value in the 2nd column**
- ↳ "false" means, it will **only search for exact matches**
- ↳ Optional: Locking the column! \$A\$2 or A\$2

Creating a Chart

- ↳ we learned that already
- ↳ mark the column "Prizes"
- ↳ open the tab "Insert"
- ↳ click on "Chart"

The CONCATENATE Function

- ↳ can put to strings of to cells into one together
- ↳ we will use it to put two datasets together
- ↳ First: We can add Spaces between the united strings
- ↳ =CONCATENATE(Cell, " ", Cell)
- ↳ the " " can make us put sth in between
- ↳ it could also be " _ "

"Data Cleaning Checklist"

Determine the size of the dataset:

Large datasets may have more data quality issues and take longer to process. This may impact your choice of data cleaning techniques and how much time to allocate to the project.

Determine the number of categories or labels:

By understanding the number and nature of categories and labels in a dataset, you can better understand the diversity of the dataset. This understanding also helps inform data merging and migration strategies.

Identify missing data:

Recognizing missing data helps you understand data quality so you can take appropriate steps to remediate the problem. Data integrity is important for accurate and unbiased analysis.

Identify unformatted data:

Identifying improperly or inconsistently formatted data helps analysts ensure data uniformity. This is essential for accurate analysis and visualization.

Explore the different data types:

Understanding the types of data in your dataset (for instance, numerical, categorical, text) helps you select appropriate cleaning methods and apply relevant data analysis techniques.

Glossary

Freitag, 11. Juli 2025 10:29

Terms and definitions for Course 4, Module 2

Clean data: Data that is complete, correct, and relevant to the problem being solved

Compatibility: How well two or more datasets are able to work together

CONCATENATE: A spreadsheet function that joins together two or more text strings

Conditional formatting: A spreadsheet tool that changes how cells appear when values meet specific conditions

Data engineer: A professional who transforms data into a useful format for analysis and gives it a reliable infrastructure

Data mapping: The process of matching fields from one data source to another

Data merging: The process of combining two or more datasets into a single dataset

Data validation: A tool for checking the accuracy and quality of data

Data warehousing specialist: A professional who develops processes and procedures to effectively store and organize data

Delimiter: A character that indicates the beginning or end of a data item

Dirty data: Data that is incomplete, incorrect, or irrelevant to the problem to be solved

Duplicate data: Any record that inadvertently shares data with another record

Field length: A tool for determining how many characters can be keyed into a spreadsheet field

Incomplete data: Data that is missing important fields

Inconsistent data: Data that uses different formats to represent the same thing

Incorrect/inaccurate data: Data that is complete but inaccurate

LEFT: A function that returns a set number of characters from the left side of a text string

LEN: A function that returns the length of a text string by counting the number of characters it contains

Length: The number of characters in a text string

Merger: An agreement that unites two organizations into a single new one

MID: A function that returns a segment from the middle of a text string

Null: An indication that a value does not exist in a dataset

Outdated data: Any data that has been superseded by newer and more accurate information

Remove duplicates: A spreadsheet tool that automatically searches for and eliminates duplicate entries from a spreadsheet

Split: A function that divides text around a specified character and puts each fragment into a new, separate cell

Substring: A smaller subset of a text string

Text string: A group of characters within a cell, most often composed of letters

TRIM: A function that removes leading, trailing, and repeated spaces in data

Unique: A value that can't have a duplicate

Datei auswählen
Keine ausgewählt