

Database Features & Components

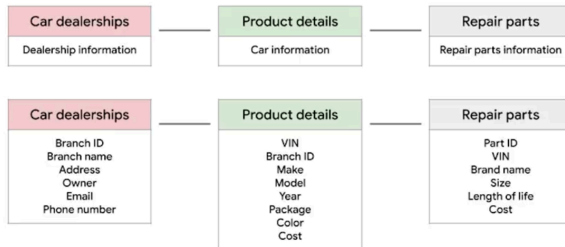
Freitag, 11. Juli 2025 10:29

Metadata

- ↳ data about data
- ↳ a data analyst thinking about how he's analysing data

Basic Database Structure Example

+ this example is called a relational database



Relational Database

- ↳ a database that contains a series of related tables that can be connected via their relationships
- ↳ one of the same fields must exist inside both tables
- ↳ in our example it's the field "Branch ID"

"In a non-relational table, you will find all of the possible variables you might be interested in analysing all grouped together. This can make it really hard to sort through. Relational databases simplify a lot the analysis process & make it easier to find data & use it across an entire database."

Normalisation

- ↳ the process of organising data in a relational database
- ↳ eliminates data redundancy
- ↳ increases data integrity
- ↳ reduces complexity in a database

Primary Key

- ↳ identifier that references a column in which each value is unique
- ↳ think of it as a unique identifier for each row in the table
- ↳ in our example "Branch ID" is the primary key
- ↳ for the "Product Details Table", "VIN" is our primary key
- ↳ if we use primary keys, they should be unique
- ↳ there can only be one primary key

Key Features

- ↳ used to ensure data in a specific column is unique
- ↳ uniquely identifies a record in a relational database table
- ↳ only one primary key is allowed in a table

Foreign key

- ↳ a field within a table that is a primary key in another table
- ↳ how one table can be connected to another table
- ↳ in our example it's "Part ID"
- ↳ each row in our "Repair Parts Table" represents one unique part
- ↳ all the other keys are foreign keys, allowing other tables to connect
- ↳ there can be multiple foreign keys

Key Features

- ↳ a column or a group of columns in a relational database, that provides a link between the data in two tables
- ↳ refers to the field in a table that's the primary key of another table
- ↳ more than one foreign key is allowed to exist in a table

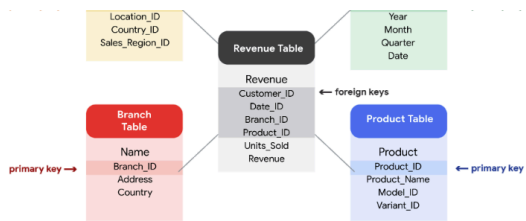
Composite Key

- ↳ a primary key constructed using multiple columns in a table

Some tables **don't require** a primary key - like a revenue table



Datei a



SQL = Game-Changer, when working with multiple databases

Metadata Management

Freitag, 11. Juli 2025 10:29

Metadata

- ↳ data about the data
- ↳ used in database management
- ↳ helps data analysts interpret contents of data within database

First Example

- ↳ check the extra-information inside a photo from our gallery
- ↳ every time we send a message - and check the message details
- ↳ it might include time, location, and more

Spreadsheets also contain all kinds of Metadata

- ↳ file type, file size, owner, last modified date, download permissions

Metadata Repositories

- ↳ specialised databases created to store and manage Metadata
- ↳ to bring together multiple sources for data analysis
- ↳ especially important, when using third-party data

Data Governance

- ↳ a process to ensure the formal management of a company's data assets

3 Types of Metadata

Descriptive

- ↳ describes a piece of data
- ↳ can be used to identify it at a later point of time
- ↳ e.g. the code on a spine of a book, called ISBN

Structural

- ↳ indicates how a piece of data is organised
- ↳ whether it is part of one, or more than one, data collection
- ↳ e.g. how the pages of the book are put together into chapters
- ↳ also keeps track of the relationship between two things

Administrative

- ↳ indicates the technical source of a digital asset
- ↳ e.g. the picture - file-type, file-size

The benefits of Metadata

Reliability

- ↳ helps us to confirm the reliability of the data we use
- ↳ make sure it's accurate, precise, relevant, timely

Consistency

- ↳ having survey data for two different sources
- ↳ using metadata to make sure, same collections methods were applied
- ↳ make sure it's organised, classified, stored, accessed

Datei auswählen
Keine ausgewählt

Access Data Sources

Freitag, 11. Juli 2025 10:29

Internal Data

↳ often called primary data

External Data

↳ often called secondary data

Comma Separated Values Files (.csv)

- + plain text files with organised table structure
- + including rows and columns
- + values in each row are separated by commas
- + widespread compatibility (*for import/export*)

We know how to do normal importing of .csv files.

Dynamic Import

IMPORTANCE function

↳ enables range of cells in one spreadsheet to be duplicated in another

↳ =IMPORTANCE ("URL", "sheet_name!cell_range")

↳ allows to automatically update spreadsheets

IMPORTHTML Function

↳ importing HTML tables or lists on a web page

↳ this process is often called "*scraping*"

IMPORTDATA Function

↳ sometimes data displayed on the web is in another form

↳ form of comma- or tab-delimited file

↳ besides this function is similar to IMPORTANCE function

Datei auswählen
Keine ausgewählt

Public Data Sets

Freitag, 11. Juli 2025 10:29

Open data helps create a lot of public datasets.
You can access them to make data-driven decisions.
Here are some resources you can use to start:

[Google Cloud Public Datasets](#)

Allow data analysts access to high-demand public datasets, and make it easy to uncover insights in the cloud.

[Dataset Search](#)

Can help you find available datasets online with keyword searches.

[Kaggle](#)

Has an Open Data search function; help find datasets to practice with

[Big Query](#)

Hosts 150+ public datasets, access and use, public health datasets

[Global Health Observatory data](#)

You can search for datasets from this page or explore featured data collections from the World Health Organization.

[The Cancer Imaging Archive \(TCIA\) dataset](#)

Just like the earlier dataset, this data is hosted by the Google Cloud Public Datasets and can be uploaded to Big Query.

[1000 Genomes](#)

This is another dataset from the Google Cloud Public resources that can be uploaded to Big Query.

[National Climatic Data Centre](#)

The NCDC Quick Links page has a selection of datasets you can explore.

[NOAA Public Dataset Gallery](#)

The NOAA Public Dataset Gallery contains a searchable collection of public datasets.

[UNICEF State of the World's Children](#)

This dataset from UNICEF includes a collection of tables that can be downloaded.

[CPS Labor Force Statistics](#)

This page contains links to several available datasets that you can explore.

[The Stanford Open Policing Project](#)

This dataset can be downloaded as a .csv file for your own use.

Keine ausgewählt

Sort & Filter Data

Freitag, 11. Juli 2025 10:29

Sorting Data

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

Freezing a Row

- ↳ select a row
- ↳ go to the tab "view"
- ↳ press on "freeze"
- ↳ select "1 row"

Multiple Sorting

- ↳ select whole sheet
- ↳ go to the tab "data"
- ↳ then "sort sheet"
- ↳ then "advanced"
- ↳ you can add several parameter
- ↳ like state and city sorting together

Filtering

- ↳ showing only the data that meets a specific criteria
- ↳ while hiding the rest
- ↳ go to "Data" and "Filter"
- ↳ you can give a word from the sheet
- ↳ and everything else will be hidden

Convert Data (*Text -> Numbers*)

- ↳ select a column
- ↳ go to "edit" - "find and replace"
- ↳ remember to click on "Match entire cell contents"
- ↳ you can type in a word and replace by a number
- ↳ or a number by a text, probably more rare

Keine ausgewählt

SQL - Big Query

Freitag, 11. Juli 2025 10:29

Big Query

- + a data warehouse
- + it belongs to Google
- + free version allows me to have 12 projects

- + you can access public data sets
- + you have a console to access (*query*) data

- + search for a data set in the **Explorer**
- + once you found something
 - ↳ use the three dots next to it
 - ↳ and press query

Basic Functions

```
SELECT *
FROM 'the_name_of_the_database'
WHERE column_name='word'
```

Name of the Database

Composed of..
Database_Name . Database_Table

If we wouldn't write `SELECT *`
↳ we would have to write every field name

Be careful about the format

↳ using ` or , or ; or >=

Good Example

```
SELECT
  duration,
  start_station_name
FROM
  `bigquery-public-data.london_bicycles.cycle_hire`
WHERE
  duration >= 1200;
```

We can also put function into SQL

↳ like `SUM()`

Sometimes you need " " to avoid errors:

↳ `favorite_food = "Shepherd's pie"`
↳ because there is a ' in the text

Use -- to create single-line notes for yourself!

You can also use /* and */ for multi-line notes!

↳ so that you remember what you were searching
↳ later in excessive usage it will be important

You can use AS to give your columns new names!

Names should never have Spaces between them!

Use CamelCase or Snake_Case

Your Query shouldn't have more than 100 characters!

↳ even through indentation is not a must (*using tabs & new lines*)
↳ still use them to make it easier for others to read your code!
↳ it's best practice!

Datei auswählen
Keine ausgewählt

SQL - Regex

Freitag, 11. Juli 2025 10:29

Regex = Regular Expressions

Defined as flexible and precise pattern-matching tools

Cheat Sheet

https://jdhao.github.io/2019/02/28/sublime_text_regex_cheat_sheet/

Examples

-- Namen, die mit "A" beginnen

```
SELECT * FROM users
WHERE REGEXP_CONTAINS(name, r'^A');
```

-- Zeigt nur valide E-Mails (vereinfachtes Muster)

```
SELECT * FROM users
WHERE REGEXP_CONTAINS(email, r'^[w._%+-]+@[A-Za-z0-9.-]+\.[A-Za-z]{2,}$');
```

-- Nur Telefonnummern mit exakt 10 Ziffern

```
SELECT * FROM contacts
WHERE REGEXP_CONTAINS(phone, r'^\d{10}$');
```

-- Sonderzeichen aus Telefonnummern entfernen

```
SELECT REGEXP_REPLACE(phone, r'^[0-9]', '') AS cleaned_phone
FROM contacts;
```

-- PLZ (Postleitzahlen) mit 5 Ziffern extrahieren

```
SELECT REGEXP_EXTRACT(address, r'\b\d{5}\b') AS zip_code
FROM customers;
```

-- Zeilen finden, die eine Zahl enthalten

```
SELECT * FROM products
WHERE REGEXP_CONTAINS(description, r'\d');
```

-- Nur URLs mit "https"

```
SELECT * FROM links
WHERE REGEXP_CONTAINS(url, r'^https://');
```

-- Alle Strings mit Bindestrich oder Unterstrich

```
SELECT * FROM strings
WHERE REGEXP_CONTAINS(col, r'[-_]');
```

Keine ausgewählt

SQL - Create Dataset

Freitag, 11. Juli 2025 10:29

Actually it is not that hard

- ↳ Create a Dataset
- ↳ Dataset can be in existing project (*driven-tenure-465609-q2*)
- ↳ Then go into the dataset
- ↳ press "create table"
- ↳ go for "upload"

Interesting is the point "Schema"

- ↳ the CSV files don't have column names
- ↳ information is like "Hannah, F, 29329"
- ↳ name, gender, count, is missing

So into schema we enter following text:

```
name:string,  
gender:string,  
count:integer
```

New Function - ORDER BY

- ↳ can make the columns be sorted
- ↳ e.g. ORDER BY count DESC
- ↳ means - order by count descending

New Function - LIMIT

- ↳ can limit the number of appearing rows
- ↳ e.g. LIMIT 5

Example with 2 new functions

```
SELECT  
*  
FROM  
`driven-tenure-465609-q2.Babynames.names_2014`  
WHERE  
gender = 'M'  
ORDER BY  
column DESC  
LIMIT  
5
```

New Select-Option

```
SELECT  
AVG(tree_dbh)
```

Datei auswählen
Keine ausgewählt

Glossary

Freitag, 11. Juli 2025 10:29

Terms and definitions for Course 3, Module 3

Administrative metadata: Metadata that indicates the technical source of a digital asset

CSV (comma-separated values) file: A delimited text file that uses a comma to separate values

Data governance: A process for ensuring the formal management of a company's data assets

Descriptive metadata: Metadata that describes a piece of data and can be used to identify it at a later point in time

Foreign key: A field within a database table that is a primary key in another table (Refer to primary key)

FROM: The section of a query that indicates where the selected data comes from

Geolocation: The geographical location of a person or device by means of digital information

Metadata: Data about data

Metadata repository: A database created to store metadata

Naming conventions: Consistent guidelines that describe the content, creation date, and version of a file in its name

Normalized database: A database in which only related data is stored in each table

Notebook: An interactive, editable programming environment for creating data reports and showcasing data skills

Primary key: An identifier in a database that references a column in which each value is unique (Refer to foreign key)

Redundancy: When the same piece of data is stored in two or more places

Schema: A way of describing how something, such as data, is organized

SELECT: The section of a query that indicates the subset of a dataset

Structural metadata: Metadata that indicates how a piece of data is organized and whether it is part of one or more than one data collection

WHERE: The section of a query that specifies criteria that the requested data must meet

World Health Organization: An organization whose primary role is to direct and coordinate international health within the United Nations system

Datei auswählen
Keine ausgewählt