

Operational Databases are Summarized and Stored

Analytical Databases are Olap online

Olap access is by business intelligence tools

Unstructured Data is raw <sup>state</sup> data

Structured data is data formatted for storage

Semistructured data - data that has been processed to a certain extent, has relevant meaning

data is accessed with SQL commands.

Record is one row of data

a File is a collection of related records

Data redundancy is bad and can cause poor security and data inconsistency

Uncontrolled data redundancy is a big issue

Data anomaly is inconsistent data, usually caused by uncontrolled data redundancy

Data Model - ERD, First step in designing a database / database blueprint

Attribute - column / field "Customer name"

One to many - Painter creates many paintings

Many to Many - Employee learns skills and skills can be learned by anyone

One to one - Store manager manages one store

Constraints - placed on data to ensure data integrity (a rule)

Schema - all of the objects grouped together

Subschema - part of the schema seen by the program

DML commands Select, insert, update, delete

Commit and Rollback

DDL Data Definition Language - Create, alter, drop

Each row in a relation is called a tuple

RDBMS - Relational Database Management system

UML - Unified Modeling language

ADT - Abstract data type

Primary Keys

Simple - Consists of one column

Composite - Consists of multiple columns

(First name + Last name + Email) (Composite Key)

Natural  
✓

Compound - multiple columns that are keys themselves

## Basic Database operations

Create

Read

Update

Delete

## 3 Types of List Anomalies

Deletion

update

Insertion

	Attribute	Attribute	Attribute
Entity			
Entity			
Entity			
Entity			

SQL - Structured query Language

ETL - Extract, transform, Load

## 4 Components of a database system are:

- Users

- Database Applications

- Database management system (DBMS)

- Database

- A Database is a self-describing collection of related records

- Tables within a database are related

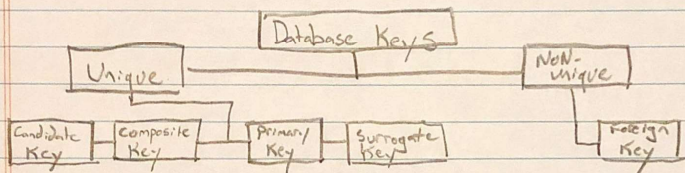
The DBMS (Database management system) - serves as an intermediary between database applications and the database

A relation is a two-dimensional table

columns are attributes

intersection between row and column are called cells

No two rows within a table can be identical



- A Composite Key is a key that is composed of two or more attributes
- A candidate key is called "candidate" because it has the potential to become the primary key. It is a unique key
- A Primary key is a candidate key chosen to be the main key for the relation
- A surrogate key is a unique, numeric value that is added to a relation to serve as the primary key. They have no meaning to users
- A Foreign key is a primary key from one table that is placed into another table
- A Null value means that no data exists
- A Functional Dependency is a relationship between attributes in which one attribute determines the value of another attribute in the same table

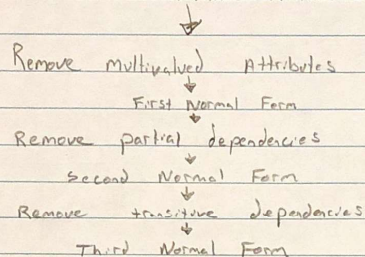
## Data Normalization

Video 2 of 8 around 35 mins in

A process of analysing a relation to ensure that it is well formed i.e. absent of the anomaly (Deletion, update or insertion)

As a general rule, a well-formed relation will not encompass more than one business concept.

Table with multivalued Attributes



## SQL Data Definition

Create - creates database objects

Alter - modify the structure and/or characteristics of existing database objects

Drop - to delete existing database objects

Count - Counts the number of rows that match the specified criteria

Min - Finds the minimum value for a specific column matching specified criteria

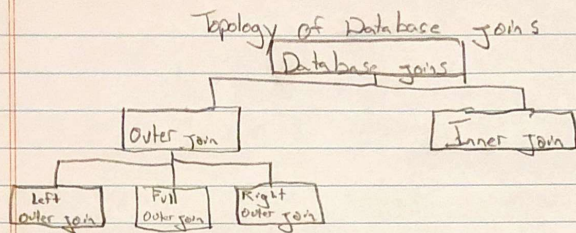
Max - Finds the maximum value for a specific column matching specified criteria

SUM - Calculates the sum (total) for a specific column for those rows matching the criteria

AVG - Calculates the numerical average (mean) of a specific column for those rows matching the criteria

STDEV - Calculates the standard deviation of the values in a numeric column whose rows match the criteria





one-to-one — ||

one-to-many — |←

zero-to-one — 0+

zero-to-many — 0←

- Functional dependency - The relationship (within the relation) that describes how the value of one attribute may be used to find the value of another attribute
- Determinant - An Attribute that can be used to find the value of another attribute in the relation

The Three most critical database administration functions are:

- 1) Concurrency control
- 2) Security
- 3) Backup and Recovery

- Concurrency - people or applications can access the same information at the same time

Dirty read - Transaction reads a modified record that has not yet been committed to the database

Inconsistent read - Transaction re-reads a dataset and finds that the data has changed

Phantom read - Transaction re-reads a dataset and finds a new record has been added

Implicit locks are issued automatically by the DBMS based on an activity

Explicit locks are issued by users requesting exclusive rights to specified data

- Table locks
- Row locks
- Column locks
- Cell locks

Dead lock - when two processes lock resources at the same time

Consistent transactions are often referred to by the acronym ACID

Atomic

Consistent

Isolated

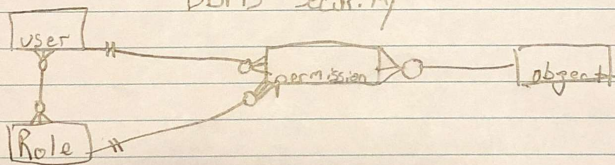
Durable

A cursor is a pointer into a set of rows that is the result set from a SQL select statement

Processing Rights define:

- who is permitted to perform certain actions
- when certain actions are allowed to be performed

### DBMS Security



BI systems fall into two broad categories:

- Reporting systems that sort, filter, group and make elementary calculations on operational data.
- Data mining applications that perform sophisticated analyses on data, analyses that usually involve complex statistical and mathematical processing.

### Operational Database

- Used for structured transaction data processing
- Current data are used
- Data are inserted, updated and deleted by users

### Dimensional Database

- Used for unstructured analytical data processing
- Current and historical data are used
- Data are loaded and updated systematically, not by users

### Data Mining Techniques

- Cluster analysis - Identifies groups of entities that have similar characteristics.
- Decision tree analysis - Classifies entities into groups based on past history.
- Regression - Produces mathematical equations that can be used to predict future events based on past observations.
- Neural Networks - Use training data to learn how to create accurate predictions / estimations.
- Market Basket Analysis - Determines patterns of associated buying behavior.



- ✓ Most data that can be encountered are best classified as semistructured
- ✓ To reveal meaning, information requires context
- ✓ The organization of data within folders in a manual file system is determined by its expected use
- ✓ An XML database supports the storage and management of semistructured XML Data
- ✓ A workgroup database supports a relatively small number of users or a specific department within an organization.
- ✓ Knowledge is the body of information and facts about a specific subject.
- ✓ Data constitutes the building blocks of information
- ✓ Unstructured data exist in the format in which it was collected
- ✓ Data inconsistency exists when different and conflicting versions of the same data appear in different places
- ✓ A data dictionary contains at least all of the attribute names and characteristics for each table in the system
- ✓ Information is the result of processing raw data to reveal its meaning
- ✓ Performance tuning relates to activities that make a database operate more efficiently in terms of storage and access speed
- ✓ XML, Extensible Markup Language is a special language used to represent and manipulate data elements in a textual format
- ✓ An application might be written by a programmer or it might be created through a DBMS utility program.
- ✓ Data is said to be verifiable if the data always yields consistent results.
- ✓ Accurate, relevant, and timely information is the key to good decision making.

One disadvantage of a database system over previous data management approaches is increased cost

End-User data is raw facts of interest to the end user

Analytical Databases focus primarily on storing data used to generate information required to make strategic decisions

A query is a specific request issued to the DBMS for data manipulation

The term islands of information refers to scattered locations storing the same basic data.

Metadata provide a description of the data characteristics and the set of relationships that link the data found within a database

An ad hoc query is a spur-of-the-moment question.

A file is a collection of related records

#### DDL Commands

- Drop
- Alter
- Rename
- Create
- Truncate

Termary Relationship involves three different entity types

Normalization - the process of organizing the fields and tables of a relational database to minimize redundancy and dependency

#### DML Commands

- Update
- Delete
- Insert
- Merge
- Select

Weak Entity is also called a dependent entity

Aggregation - is used for populating summaries that can be used at the staging area of ETL

Data warehouse

- Centralized repository of information
- Organized around relevant subject areas
- Provides platform for queries
- Used for analysis and not transactional processing
- Data is non-volatile
- Target location for integrating data from multiple sources

Every non-Key must provide a fact about the Key, the whole Key, and nothing but the Key

- "the Key" Ensures 1NF
- "the whole Key" Ensures 2NF
- "Nothing but the Key" Ensures 3NF