

# ANALYZING BIG DATA - I

---

Week 2 – Introduction to (Relational) Data Bases

## Agenda

- **Defining Data bases (DB's):** History of DB's, Types of DB's, (Dis)Advantages of DB's
- **Understanding** relational DB's:
  - Elements of a relational-DB schema, tables, attributes, relations, primary keys, foreign keys
  - Basic modeling exercises
- **Interacting** with relational DB's: One will always need
  - A DBMS: We use MySQL workbench
  - A Programming language that is SQL – based: We use MySQL
    - History of MySQL
    - SQL vs MySQL vs PostgreSQL vs Teradata SQL
  - We illustrate everything with: Geography DB
    - Create: Schemas, Tables, Relationships, Attributes, Add records to a table
    - Drop: Tables, DB's
    - Very Basic Searches

# DATABASES

---

## Databases – History

- Storing and manipulating data has **always** been a major computer application
- Early 1960's – Charles Bachman at General Electric designs the first Integrated Data Center (Store)
  - Nobel Prize in Computer Science.... Does not exist ☹
- Late 1960's –
  - IBM jumps in and develops the Information Management System (IMS)
  - The SABRE system is developed by American Airlines (IBM help). (Fly tickets)
- 1970's – Edward Codd proposes the relational model.
  - We will study the *relational data model* today for data storage
  - *Relational data model* became universal and Programming languages specialized to work with it are developed.
- 1990's – Today : The relational data model receives a new Boost. Why?

## Databases – General Description

- Database: Any set of data held in a computer in such a way that can be easily accessed, managed or updated.
  - Definition is too wide (useless).
  - Solution: Further classify.
- Types of databases:
  1. **Nonrelational DB's** (or NoSQL Based)
    - Examples
      - Key-value Stores: (Redis) and Amazon DynamoDB
      - Wide Column Stores: [Cassandra](#), [Scylla](#), and [HBase](#)
      - Graph Databases: (Neo4J, DataStax)
      - Document oriented DB's: MongoDB and Couchbase
      - Search engines: ElasticSearch
    - Advantages: Very Flexible
    - Disadvantages: Still at its infancy
  2. **Relational DB's** (or SQL Based)
    - Based on Entities and Relationships

## Databases – General Description

- Database: Any set of data held in a computer so that it can be easily accessed, managed and updated.
  - Definition is too wide (useless)
  - Solution: Further classify.
- Types of databases:
  1. **Nonrelational DB's** (or NoSQL Based)
    - Examples
      - Key-value Stores: (Redis) and Amazon DynamoDB
      - Wide Column Stores
      - Graph Databases
      - Document oriented DB's
    - Advantages: Very Flexible
    - Disadvantages: Still at its infancy
  2. **Relational DB's** (or SQL Based)
    - Based on Entities and Relationships
    - Rely on Database Management System

## Databases – General Description

- Database: Any set of data held in a managed and updated.
  - Definition is too wide (useless)
  - Solution: Further classify.

### Types of databases:

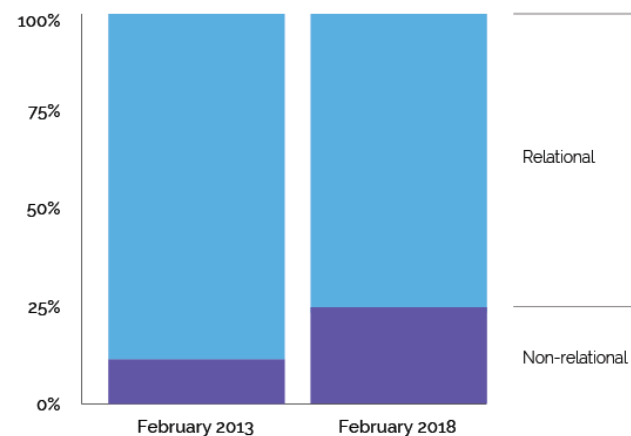
#### 1. Nonrelational DB's (or NoSQL Based)

- Examples
  - Key-value Stores: (Redis) and Amazon D
  - Wide Column Stores
  - Graph Databases
  - Document oriented DB's
- Advantages: Very Flexible
- Disadvantages: Still at its infancy

#### 2. Relational DB's (or SQL Based)

- Based on Entities and Relationships
- Rely on Database Management System

Popularity (percentage) Relational Databases vs. Non-Relational Databases



Source: [https://db-engines.com/en/ranking\\_trend](https://db-engines.com/en/ranking_trend)

# UNDERSTANDING RELATIONAL DB'S

---



## Relational DB's

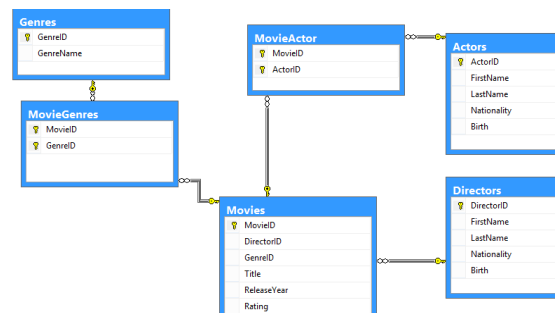
- Relational DB's: Collection of data describing the activities of an **specific** organization and it is composed by **Entities** and **Relationships**:
  - **Entities**: An **entity** is any object in the system that we want to model and store information about.
    - Movies DB: Genres, Actors, Directors, Movies
  - **Relationships**: Restrictions linking each of the entities
    - Movies DB: Movies has many actors and actors appear in many movies

## Relational DB's - Examples

See word document

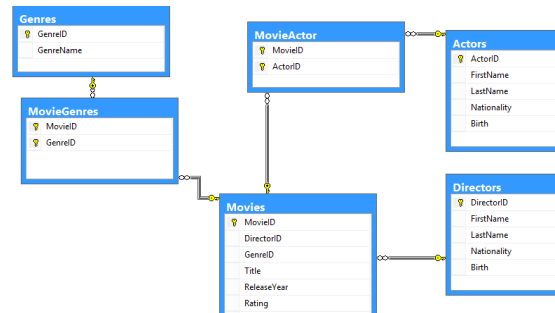
## Relational DB's - Schema

- The Key element of a Relational DB is its schema
- It is graph which **must** include the following elements
  1. Entities (tables)
    - Attributes (or columns, or variables)
    - All the Primary keys
    - All the Foreign keys
  2. Relations
    - Relations **between real entities** must indicate its type: 1-to-many, many-to-many



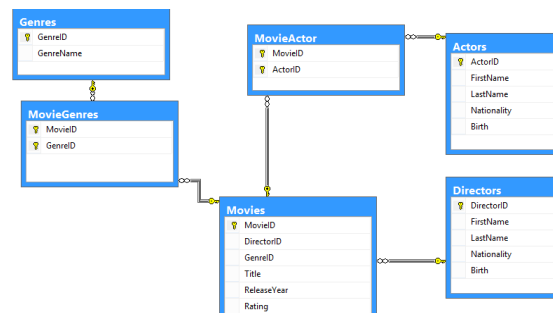
## Relational DB's - Schema - Tables

- The Key element of a Relational DB is its schema
- It is graph which **must** include the following elements
  1. Entities (tables)
    - Attributes (or columns, or variables)
    - All the Primary keys
    - All the Foreign keys

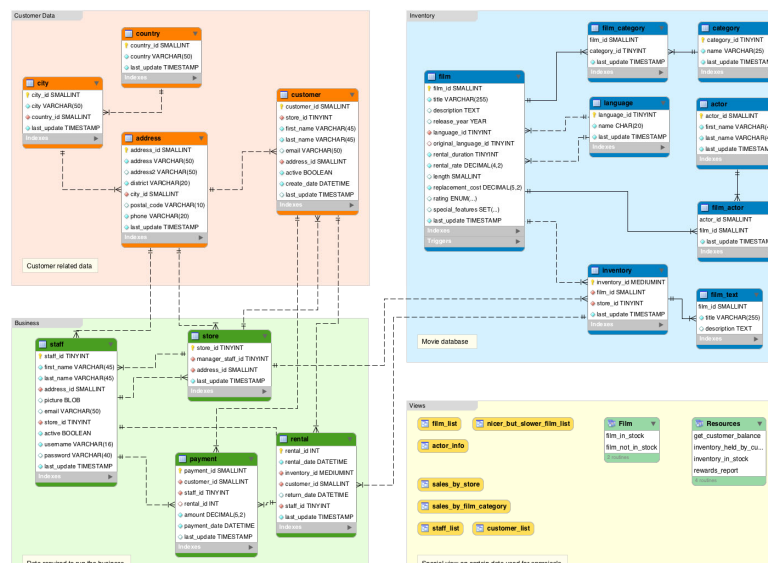


## Relational DB's - Schema - Relations

- One-to-One: *Not a Problem (trivial)*
- One-to-many: Not a Problem
- Many-to-many: Almost Not a Problem
  - > Break it into two One-to-One Relationships
- Important: An Schema is a hierarchical structure

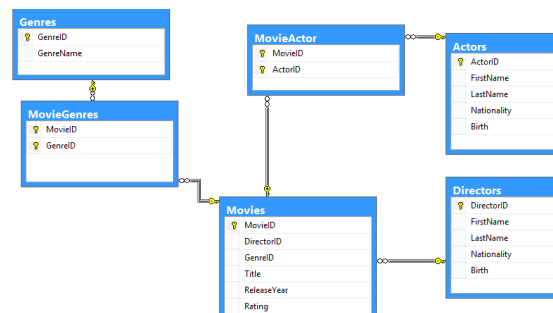


# Relational DB's - Schema - Examples



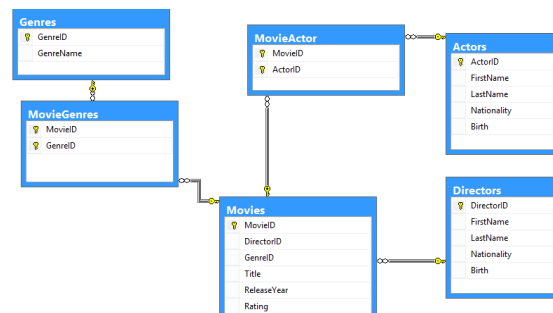
## Relational DB's - Schema: In Practice

- Data Base is a whole field in Computer Science
- You don't need to know everything !!!!
- As a data scientist **you just need to be able to:**
  - **Design** a basic schema for an organization based on a description of it.
  - **Describe** an organization based on the schema of the data base
  - **Recover an existing schema:**
    - My-SQL will do it for you
    - If your company does not use My-SQL. Still can be done
      - No worries 😊

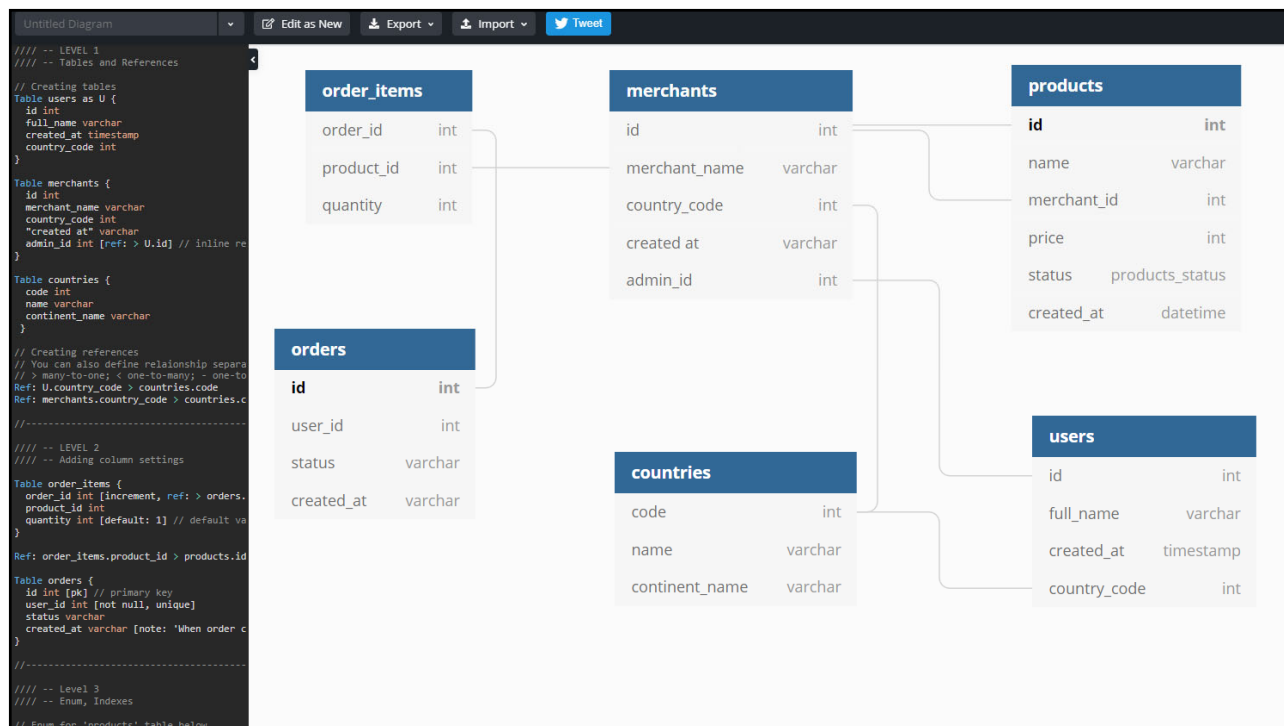


## Relational DB's - Schema - Design

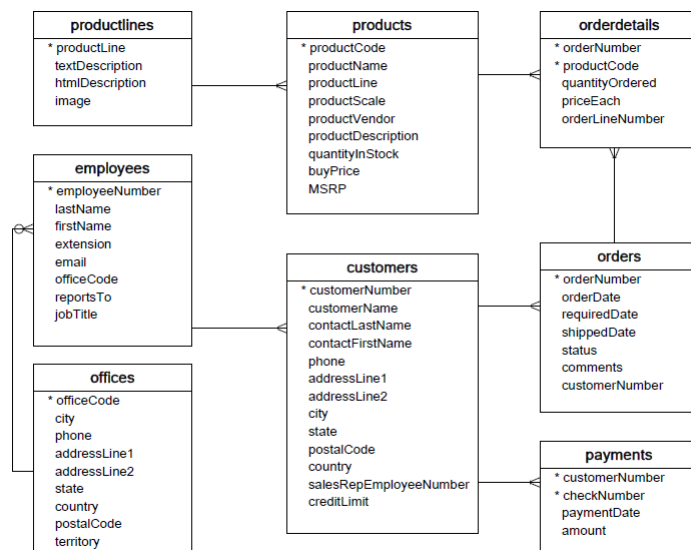
- Step one: List all the relevant entities
  - List all the attributes you want
  - Specify if they have to be stored as text or numerically
- Step two:
  - Identify as many relationships as possible between the entities.
  - Determine whether is 1 to many, 1 to 1 or many to many
    - If it is many to many -> Break it
- Go to MySQL or to <https://dbdiagram.io/home>





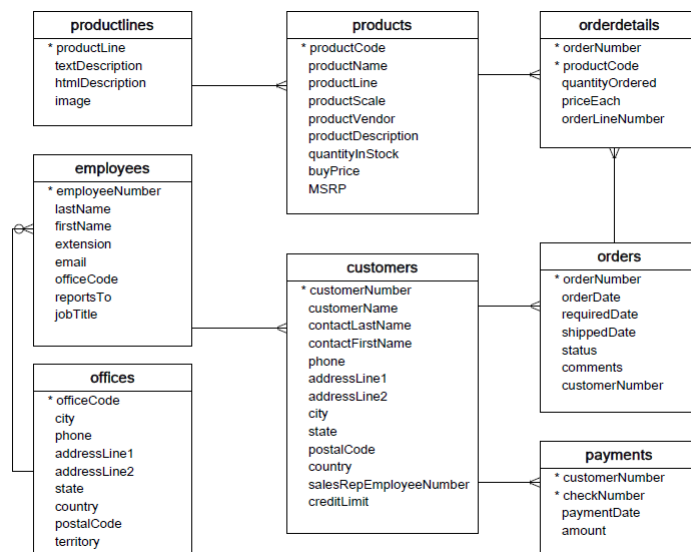


## Relational DB's – Describe an Organization based on the Schema



## Relational DB's - Recovering Schemas

- Done in My-SQL
- When you get to your first job, your boss may not give you the schema.
- They may just give you access to the database



## Relational DB's – Advantages

- Mature Technology
- Data Integrity and Security
- Efficient (and concurrent) Data Access
- Data Independence
- Data administration

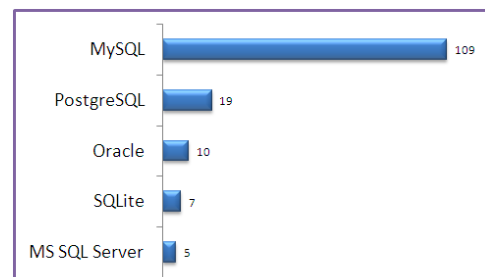
## Relational DB's – Disadvantages

- Relational databases are VERY costly to install
- DBMS are **EXTREMELY** costly to change
  - Switching costs can become a serious problem
  - DBMS create “*Path dependence*”
    - IBM
    - Oracle
    - Amazon

# INTERACTING WITH RELATIONAL DB'S

---

Welcome to *SQL - Pronounced Sequel* 😊



## SQL – History

SQL - Programming Language many dialects



## SQL – From SQL to My-SQL

- In what follows we switch to My-SQL we are going to work directly with the DBMS
- The goal is to replicate all what he have done in a hands-on approach using My-SQL

# Visiting My-SQL workbench

The screenshot shows the MySQL Workbench interface. The main editor window displays the following SQL script:

```

1  -- Preliminaries
2  * DROP DATABASE IF EXISTS db_countries;
3  * CREATE DATABASE db_countries;
4  * USE db_countries;
5
6  -- Creating the Tables
7  * CREATE TABLE continents(
8      id          INT UNSIGNED NOT NULL AUTO_INCREMENT, # Unique ID for the record
9      NAME        VARCHAR(100) NOT NULL,               # Name of the continent
10     SURFACE      INT,                                # Surface of the continent in number of s
11     POPULATION   BIGINT,                             # Number of people. Integer number
12     PRIMARY KEY (id)                                # Make the id the primary key
13 );
14
15
16 * CREATE TABLE countries(
17     ID_COUNTRY   INT UNSIGNED NOT NULL AUTO_INCREMENT, # Unique ID for the record
18     NAME         VARCHAR(100) NOT NULL,               # NAME of the country
19     CONTINENT    INT UNSIGNED NOT NULL,              # NAME of the continent where the country is located
20     SURFACE      INT,                                #
21     POPULATION   INT,                                #
22     CAPITAL      INT,                                #
23     PRIMARY KEY (ID_COUNTRY)                         #
24 );
25
26
27 * CREATE TABLE cities(
28     ID_CITY      INT UNSIGNED NOT NULL AUTO_INCREMENT, #
29     CITY_NAME_CHAR VARCHAR(150) NOT NULL,             #
30     SURFACE      VARCHAR(150) NOT NULL,               #
31     POPULATION   VARCHAR(150) NOT NULL,               #
32     PRIMARY KEY (ID_CITY)                            #
33 );
34
35

```

The interface includes a sidebar on the left with a 'Schemas' tree showing 'db\_countries' and its tables. The bottom panel shows the 'Output' tab with a message: '1 14:02:52 USE db\_countries 0 rows affected'.

## Creating a DB

- On My-Sql workbench

## Creating a Table

- On My-Sql workbench

## Attribute type - NUMERIC

- Integers
  - TINYINT/ BOOL/BOOLEAN
  - SMALLINT
  - MEDIUMINT
  - INT/INTEGER
  - BIGINT
- Floating point
  - FLOAT: single precision 7 digits
  - DOUBLE: double precision 15 digits.
- Fixed Point
  - DECIMAL(6,2) represents 6 digits of which 2 will be decimal points.

## Attribute type - STRINGS

- Character Strings:
  - Fixed length: CHAR (for example CHAR(8)).
  - Variable length: VARCHAR (p.e. VARCHAR(255)).
  - Very long text: TEXT
- Un-ordered Sets: ENUM
  - ENUM('Dean','Full-Professor','Associated-Professor','Assistant-Professor')
- Ordered Sets: SET
  - SET('red','green','yellow')

## Attribute type - DATE

- DATE: just the date
- DATETIME: Calendar time with hours and minutes and
- YEAR: 4 digit year
- YEAR(2): 2 digit year (don't use it)
- TIMESTAMP:

## Altering Tables – Add/Remove Column

- On My-Sql workbench



## Altering Tables – Add/Remove Foreign Key

- On My-Sql workbench

## Reverse Engineering an Schema

- On My-Sql workbench

## Adding a Column to Table

- On My-Sql workbench

## Reverse Engineering an Schema

- On My-Sql workbench

## Introducing HW 2

