

## SQL

Fetch data and process it to derive valuable insights for decision making in organization

### Week 1

#### DB admin and Data scientist

- Compare and contrast roles of DB admin and Data scientist
- Data scientist primarily use SQL for data retrieval, Combine and create tables, write queries for analysis and build test models

#### Database Administrator or Data Scientist

##### Database Administrator

Manages/governs entire database  
Gives permissions to users  
Determines access to data  
Manages and creates tables  
Uses SQL to query and retrieve data

##### Data Scientist

End user of a database  
Uses SQL to query and retrieve data

## SQL

- SQL – Insert Query Update Modify [Create Read Update Delete : operations of database]
- Non procedural [cannot write applications]
- Descriptive statements to interact with database
- DBMS have dialect – SQL can translate

## Data models

- Think about problem to solve before writing query
- Understand business process for which data is modeled, organized and structured
- Speeds up coding, less rework, improve accuracy
- RDB enables to write queries against structured data where relations between all tables are clearly defined

### Types of models

- Relational model: Scalable and Structured [modeled] data for manipulation, analysis and logical querying using keys
- Transactional model: Operational data for transactions

1960	Hierarchical	Difficult to represent M:N relationships (Hierarchical only)
1969	Network	Structural level dependency No ad hoc queries (record-at-a-time access) Access path predefined (navigational access)
1970	Relational	Conceptual simplicity (structural independence) Provides ad hoc queries (SQL) Set-oriented access
1976	Entity Relationship	Easy to understand (more semantics) Limited to conceptual modeling (no implementation component)
1978	Semantic	More semantics in data model Support for complex objects
1985	Object-Oriented	Inheritance (class hierarchy) Behavior Unstructured data (XML)
1990	Extended Relational (O/R DBMS)	XML data exchanges
2009 Big Data	NoSql	Addresses Big Data problem Less semantics in data model Based on schema-less key-value data model Best suited for large sparse data stores

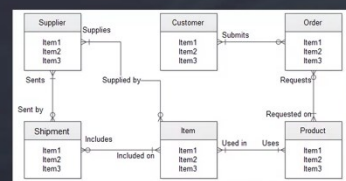
## Modeling

- Show relationships[one to many, many to many]
- Show links[Primary – uniquely ID rows & foreign keys – group of columns to ID rows]

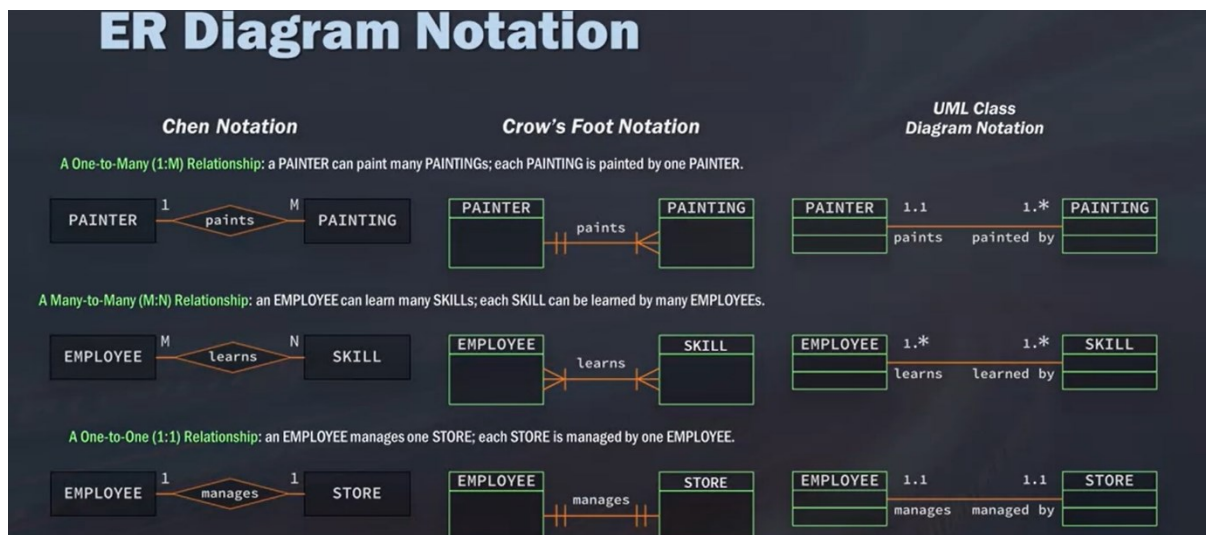
## ER Diagrams

### ER model

Is composed of entity types and specifies relationships that can exist between instances of those entity types



## Notations in ER diagram



## Reading material

- [What is SQL and How is it Used?](#)
- [NTC Hosting: Structured Query Language \(it's worth exploring this site, not just this singular posting\)](#)
- [SQLite Tutorial](#)
- [Norwalk Aberdeen: Entity-Relationship Diagrams \(9 Minute YouTube Video\)](#)
- [Star Schema vs. Snowflake Schema](#)
- [Explain Star Schema & Snow Flake Design \(5 Minute YouTube Video\)](#)
- [Data Modeling 101](#)
- [What is Data Modeling - An Introduction for Business Analysts](#)
- [Wikipedia: Data Modeling](#)
- [Dataconomy: SQL vs. NoSQL - What You Need to Know](#)
- [TechRepublic: NoSQL keeps rising, but relational databases still dominate big data](#)
- [SiliconRepublic: Data Science Skills: Is NoSQL Better than SQL?](#)

## Week 2

### SQL Statements

- **NULL** : absence of everything
- **Empty string** : 0, spaces are present
- **Comments** : block comment /\* \*/ , Line comment --

#### Wildcard % \_

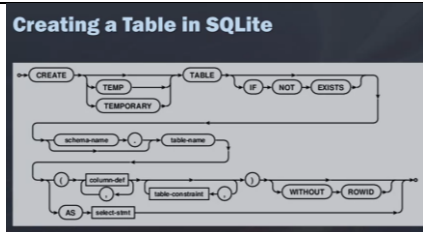
- Used with LIKE for string parsing, Does not accept NULL
- All wildcards do not work with all DBs eg : \_ DB2, [] sqlite
- Helpful but reduces query performance; use operators as much as possible
- Add % sign in "optional" positions
- %pizza → ending with pizza
- %pizza% → has pizza as part of string
- %@%.com → match any email id
- Pi\_\_a → any 2 letters accepted between Pi and a

#### Arithmetic operations () ^ / \* + -

**SELECT a, (b-c)\*d as e FROM table1**

#### Typical precedence

- **SELECT** <column>
- **FROM** <table>
- **WHERE** <row level filter>
- **GROUP BY** <group by column>
- **HAVING** <condition on groups>
- **ORDER BY** <column with order>

Statement	Syntax	Description
<b>CRUD operations</b>		
SELECT – FROM	SELECT <columns> FROM <table> LIMIT <first N records>	
CREATE TABLE	CREATE TABLE < table > ( <col1> CHAR(10) PRIMARY KEY <col2> INT(10) NOT NULL . . )	
CREATE TEMPORARY TABLE	CREATE TEMPORARY TABLE <temp_table> ( SELECT * FROM <table> WHERE <col1> = [match value] )	Create fast temporary table, deleted after session terminated simplify queries by creating a subset, and then joining to that subset, and derive a new calculation.
INSERT INTO	INSERT INTO < table >( <col1>, <col2>, ... ) VALUES( 'val1', 'val2', ... )	



HAVING	SELECT <columns> FROM <table> GROUP BY <col1>, <col2>.. HAVING COUNT(*) <operator> <match value>	<b>GROUP BY – HAVING for group wise filtering</b>
<b>Distinct function</b>		
DISTINCT	SELECT COUNT(DISTINCT <column>) AS <countDistCol> FROM <table>	Without Distinct, SQL considers all data And ignore duplicates

### Reading material

#### SQL for R

- [SQLDF Package](#)
- [Documentation](#)
- [Examples](#)

#### SQL for Spark

- [Overview](#)
- [Documentation](#)

#### SQL with Hadoop

- [Hive Overview](#)
- [Documentation](#)

#### SQL for Python

- [Python-SQL Package Documentation](#)

### Week 3

#### Subqueries [SELECT FROM WHERE IN (SELECT FROM ...)]

- **Definition:** Query inside other query with additional filtering criteria; Work on single and multiple tables  
**E.g. Know region of each customer who had order with freight more than 100**  
SELECT CustomerID, Region FROM Customers  
WHERE CustomerID in (  
    SELECT CustomerID FROM Orders  
    WHERE Freight > 100  
);
- **Calculated fields**  
**E.g. Total order placed by every Customer**  
SELECT Customer\_name, States,  
(  
    SELECT COUNT(\*) AS Orders FROM Orders  
    WHERE Orders.CustomerID = Customers.CustomerID  
)  
FROM Customers  
ORDER BY Customer\_name;
- **Limitations:** Deeply nested subqueries reduce query performance; Subquery Selects only one column at a time

#### Key Fields

- Data Segmentation i.e. Breaking data into tables logically models Business problem
- It makes data Scalable, easier for manipulation and analysis, makes storage efficient
- Logical querying using Keys which indicate attributes through which tables are related

#### Joins

- Instead of duplicating tables; tables can be **joined** using keys – a nonphysical single line query
- Associate related records from different tables on the fly

**Cross join :** It is Cartesian Join each row from T1 with each row of T2

- computationally taxing because if you have a table with just ten records in it and the second table with ten records, just performing a Cross Join is already going to increase it to 100.

**Inner Join:** matching value from both table [INTERSECTION]

- Computationally taxing to find ON conditions

**Left Join:** Values from left table with matching value from right table

**right Join:** Values from right table with matching value from left table

**Full outer Join:** Return all records where there is a match in *either* table one *or* there's a match in table two

**Union :** Combine results from multiple SELECT statements by **stacking results**

All SELECT columns must be in same sequence and have same data types

**Aliases:** table names short formed for better access

E.g. SELECT O.order\_ID FROM orders O;

**Pre-qualifiers :** subqueries for nested joins

E.g. SELECT P.product\_name, O.unit\_price, S.company

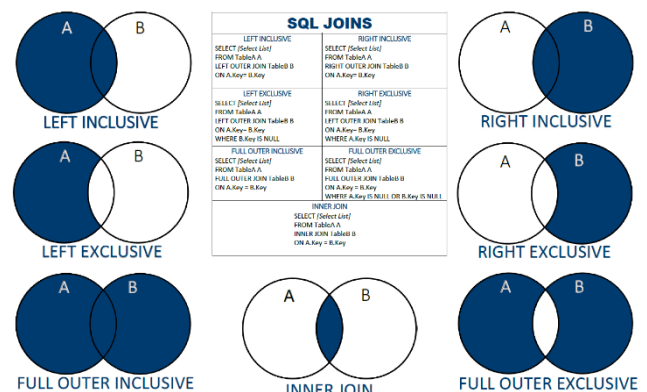
FROM (

    (suppliers S INNER JOIN products P ON S.supplierID = P.supplierID )

    INNER JOIN orders O ON P.productID = O.productID ) ;

#### Reading material

- [Thinking in SQL vs Thinking in Python](#)
- [Difference Between Union and Union All - Optimal Performance Comparison](#)



## Week4

**Concatenating** : Link columns together - Use pipes

E.g. SELECT firstname, surname, firstname || surname FROM employees → firstname, surname, firstnamesurname

**Trimming** : TRIM everything from front and back OR RTRIM, LTRIM

E.g. SELECT TRIM (" abc hdh ") AS trimmedString; → "abc hdh"

**Substring function** : select part of string - SUBSTR(<STRING>, A, B) gives substring starting at A<sup>th</sup> character and return B subsequent characters

E.g. SELECT name, SUBSTR("Jonathan", 2,3) FROM employees → ona

**Change Case** - UPPER, LOWER, UCASE, LCASE

E.g. SELECT UPPER "myname" FROM employees → MYNAME

**Datetime function** : Different databases support different datatype : DATE, TIME, DATETIME, JULIANDAY, STRFTIME

STRFTIME : extract certain elements of a date/time string

E.g.

- SELECT birthdate, STRFTIME("%y", birthdate) AS birthyear from birthdays
  - SELECT DATE('now') → return present system date
  - SELECT birthdate, DATE(("now") - birthdate) AS Present\_Age from birthdays → find current age from birth date
- Time strings are extracted from DATETIME object

**Case statements** : Used inside SELECT, INSERT, DELETE, UPDATE to Categorize or Bin data

CASE WHEN THEN ELSE → IF THEN ELSE

## Simple Case Statement

```
CASE input_expression
  WHEN when_expression THEN
    result_expression [ ...n ]
  [ ELSE else_result_expression ]
END
```

	employeeid	firstname	lastname	city	calgary
1	1	Andrew	Adams	Edmonton	Other
2	8	Laura	Callahan	Lethbridge	Other
3	2	Nancy	Edwards	Calgary	Calgary
4	5	Steve	Johnson	Calgary	Calgary
5	7	Robert	King	Lethbridge	Other
6	6	Michael	Mitchell	Calgary	Calgary
7	4	Margaret	Park	Calgary	Calgary
8	3	Jane	Peacock	Calgary	Calgary

```
SELECT
employeeid
,firstname
,lastname
,city
,CASE City
  WHEN 'Calgary' THEN 'Calgary'
ELSE 'Other'
END calgary
FROM Employees
ORDER BY LastName, FirstName;
```

## Search Case Statement

```
SELECT
trackid
,name
,bytes
,CASE
  WHEN bytes < 300000 THEN 'small'
  WHEN bytes >= 300001 AND bytes <= 500000 THEN 'medium'
  WHEN bytes >= 500001 THEN 'large'
ELSE 'Other'
END bytescategory
FROM
tracks;
```

trackid	name	bytes	bytescategory
1	Elton John's Greatest Hits	10112	small
2	Flow Sports	101200	small
3	A. Martin	211017	small
4	Ogden	211111	small
5	Commercial 1	301000	medium
6	The Road Problem	301000	medium
7	Commercial 2	301000	medium
8	Beavis	301000	medium
9	Comedy Film	301000	medium
10	Donna Estee	301000	medium
11	Matteo Ester	301000	medium
12	Michael (Michael)	301000	medium
13	John - Churchill's Speech	301000	medium
14	Urbain Act 1. Surfer (Chuchill)	301000	medium
15	Lamentations of Jeremiah (First Part) Script Lamentation	301000	medium
16	Calgary Orchestra	301000	medium
17	Frederick For No People	301000	medium
18	Democrat	301000	medium

## Views

- A view is a stored query – an illusion of table
- Add remove columns without changing schema and database write limitations
- Tryout and encapsulate complex queries without ETL

SELECT \* FROM <myView>

DROP VIEW <myView>

## Creating a View

```
CREATE VIEW my_view
AS
SELECT
r.regiondescription
,t.territorydescription
,e.Lastname
,e.Firstname
,e.Hiredate
,e.Reportsto
FROM Region r
INNER JOIN Territories t on r.regionid = t.regionid
INNER JOIN EmployeeTerritories et on t.TerritoryID = et.TerritoryID
INNER JOIN Employees e on et.employeeid = e.EmployeeID
```

regiondescription	territorydescription	lastname	firstname	hiredate	reportsto
1	Western	Davidson	Hancy	10-11-1991 12:00:00 AM	2
2	Eastern	Hewson	Hancy	5-1-1991 12:00:00 AM	2
3	Western	Steffens	Faller	1-14-1992 12:00:00 AM	2
4	Eastern	Reedford	Faller	1-14-1992 12:00:00 AM	2
5	Eastern	Georgewine	Faller	1-14-1992 12:00:00 AM	2
6	Eastern	Endine	Faller	1-14-1992 12:00:00 AM	2
7	Eastern	Cambridge	Faller	1-14-1992 12:00:00 AM	2
8	Eastern	Downey	Faller	1-14-1992 12:00:00 AM	2
9	Eastern	Louisville	Faller	1-14-1992 12:00:00 AM	2
10	Southern	Atlanta	Leavelle	4-1-1992 12:00:00 AM	2
11	Southern	Tampa	Leavelle	4-1-1992 12:00:00 AM	2
12	Southern	Atlanta	Leavelle	4-1-1992 12:00:00 AM	2
13	Southern	Tampa	Leavelle	4-1-1992 12:00:00 AM	2
14	Eastern	Richfield	Pravack	5-1-1991 12:00:00 AM	2
15	Eastern	Greenbore	Pravack	5-1-1991 12:00:00 AM	2
16	Eastern	Cary	Pravack	5-1-1991 12:00:00 AM	2
17	Eastern	Providence	Pravack	10-17-1991 12:00:00 AM	2
18	Eastern	Providence	Pravack	10-17-1991 12:00:00 AM	2
19	Eastern	Edison	Pravack	10-17-1991 12:00:00 AM	2
20	Eastern	New York	Pravack	10-17-1991 12:00:00 AM	2
21	Eastern	New York	Pravack	10-17-1991 12:00:00 AM	2
22	Eastern	Madison	Pravack	10-17-1991 12:00:00 AM	2
23	Eastern	Longport	Pravack	10-17-1991 12:00:00 AM	2
24	Western	Phoenix	Pravack	10-17-1991 12:00:00 AM	2
25	Western	Scottsdale	Pravack	10-17-1991 12:00:00 AM	2

## Data Governance and Profiling

Application of SQL in Data Science

**Data Governance**: looking at either descriptive statistics [columns MIN, MAX, DISTINCT, NULL, AVG] or different information on the data Rows [table size, last updated]

**Profiling** : Read and write policies about data site, clean up environment, escalation process

### Data understanding ↔ Business understanding

- Unspoken requirements → requirement gathering

### While profiling data:

- Details in data, create data model
  - draw out the different tables on a piece, of paper, basically creating my own data model and map.
- Consider joins and calculations necessary
- Data quality and format issues
- Test after each JOIN regressively
  - testing block by block allows you to find where any problems or issues occur
- Business rules – Date changes and Indicators
  - Has the data changed? Are the business rules different? Do you need to update and change the data indicators? Does anything need to be updated?

### Reading material

- [SQL Authority: SQL Puzzles](#)
- [SQLZOO](#)

## Beware of the Unspoken Need

“We want to predict whether or not a customer is likely to buy our product.”

Which customers?

What product?

What is/should be excluded?

What is/should be counted from past?