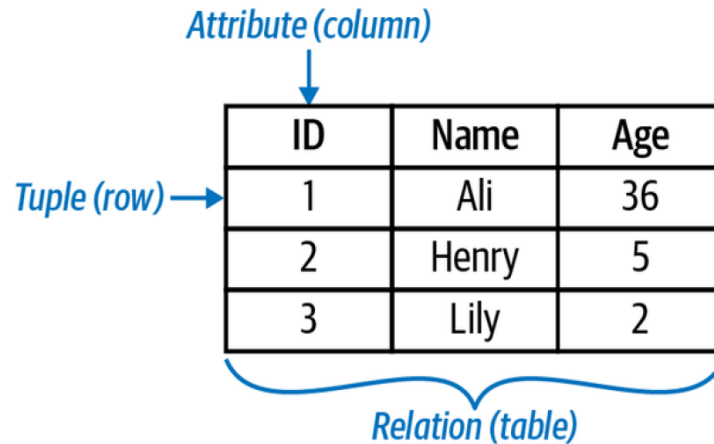


SQL Notes

SQL is a structured query language that involves using relational databases, which requires pre-defined schemas.



CRUD Operations are most used in SQL. (Create, Read, Update, Delete)

SQL statements are divided into 4 main categories:

1. Data Definition Language (DDL)
 - CREATE (Create new database/table)
 - ALTER (Modify database/table structure)
 - DROP (Delete database or table)
 - TRUNCATE (Remove table records)
2. Data Manipulation Language (DML)
 - SELECT (Retrieve data from table)
 - INSERT (Insert new records into table)
 - UPDATE (Update existing records from a table)
 - DELETE (Remove existing records from a table)
3. Data Control Language (DCL)
 - GRANT (Assign privilege to users for accessing data)
 - REVOKE (Remove privilege to users for accessing data)
4. Transaction Control Statement (TCS)
 - COMMIT (Permanent work save into database)
 - ROLLBACK (Restore database to previous form since last commit)
 - SAVEPOINT (Create save point for future rollbacks)
 - SET TRANSACTION (Setting transaction to read-write/read only access)

SQL Query Basic Template

SELECT <column_names>

FROM <table_name>

WHERE <condition on columns from table>

GROUP BY <column_names>

HAVING <condition on grouped-by columns>

ORDER BY <column_names>

SQL works in the following order of statement execution:

1. Gathers all data with FROM clause
2. Filters the data with WHERE clause
3. Groups rows together with GROUP BY clause
4. Filters grouped rows with HAVING clause
5. Specifies columns to display with SELECT clause
6. Sorts the results with ORDER BY clause

Types of RDBMS and its database tools:

1. SQLite: DB Browser for SQLite
2. MySQL: MySQL Workbench
3. Oracle: Oracle SQL Developer (PL/SQL – Procedural Language Extension to SQL)
4. PostgreSQL: pgAdmin
5. SQL Server: SQL Server Management Studio (T-SQL – Transact SQL)

ANSI (American National Standards Institute) standard in SQL refers to SQL code that will run in any RDBMS software.

Identifiers vs Aliases

Identifiers are name of database objects

Aliases rename column or table temporarily, mostly useful in subqueries.

Statements vs Clauses

Statements are blocks of code that starts with a SQL keyword like SELECT and ends with a semicolon.

Clauses are specific sections of the statement that refers to specific SQL keywords like WHERE, FROM etc.

Single vs Multi-Line Comments

Single Line Comment:

-- This is a single line comment

Multi Line Comment:

/* These are
multi line comments */

Single vs Double Quotes

Single quote: Used for string reference

Double quote: Used for identifier reference

Wildcard expressions

%: Represents any n number of characters

_: Represents any single character

Note that these wildcard expressions are used together with LIKE or NOT LIKE keyword in WHERE clause.

SQL Data Types

1. Numeric (INT, DECIMAL, FLOAT)

- INT: Used for values that do not allow for decimals (i.e. 45)
- DECIMAL: Used for fixing number of decimals (i.e. 24.524)
- FLOAT: Store limited number of decimals with power notation (i.e. $2.4524 * 10^5$)

Note: MySQL has the option of setting numeric variables as positive only using UNSIGNED keyword.

2. String (CHAR, VARCHAR, TEXT, NCHAR, NVARCHAR)

- CHAR: Stores fixed length of characters as ASCII data
- VARCHAR: Stores maximum length of specified characters as ASCII data
- TEXT: Used for storing long strings of text like paragraphs
- NCHAR: Stores fixed length of characters as Unicode data
- NVARCHAR: Stores maximum length of specified characters as Unicode data

Note: Unicode data refers to non-ASCII characters (non-English mostly)

3. Date (DATE, TIME, DATETIME, TIMESTAMP, YEAR – For MYSQL)

- DATE: YYYY-MM-DD
- DATETIME: YYYY-MM-DD hh:mm:ss
- TIMESTAMP: YYYY-MM-DD hh:mm:ss UTC
- TIME: hh:mm:ss
- YEAR: YYYY

Note: DATETIME variable type does not store time zone, while TIMESTAMP variable type does store time zone.

4. Boolean (BOOLEAN): FALSE value as 0 and TRUE value as 1

Note: Boolean variable type is currently not supported in Oracle and SQL server.

5. External files like images, documents etc.

Approach 1: Store links to files using VARCHAR variable type

Approach 2: Convert files to binary format and store files using BLOB variable type

Common SQL Operators

Logical Operators

AND: Returns TRUE if both conditions are true or otherwise

OR: Returns TRUE if either condition is true or otherwise

NOT: Returns TRUE if condition is FALSE or otherwise

Comparison Operators (Used in conditional statements)

=: Equality

!=, <>: Inequality

<: Less than

<=: Less than or equal to

>: Greater than

>=: Greater than or equal to

BETWEEN: Checks if value lies within given range (inclusive)

EXISTS: Checks if row exist in subquery (also known as semi-join)

IN: Checks if value is contained within a list using brackets symbol '()' (Note that having a single null value within the list will always result in FALSE for conditional statements. Thus, recommend using **NOT EXISTS** keyword instead.)

IS NULL: Checks if a value is null

IS NOT NULL: Checks if a value is not null

LIKE: Checks if value matches a simple pattern (Used with wildcard expressions)

Math Operators

+: Addition

-: Subtraction

***: Multiplication**

/: Division

?: Modulo (remainder)

&: Bitwise AND

|: Bitwise OR

^: Bitwise XOR

Aggregate Functions

Aggregate functions perform calculations on multiple rows to return a single result.

Aggregate functions are mostly used in SELECT and HAVING clause

COUNT(): Counts number of values for a given attribute

SUM(): Sum of values for a given attribute

AVG(): Average of values for a given attribute

MIN(): Minimum value for a given attribute

MAX(): Maximum value for a given attribute

LEAST(): Minimum value from set of attributes for every row

GREATEST(): Maximum value from set of attributes for every row

Note that most of these aggregate functions in SQL applies to only non-null values, except for COUNT() function that applies to both null and non-null values.

Numeric Functions

Numeric functions are usually applied to attributes with numerical values only.

Math functions:

ABS(x): Returns absolute value of x

SIGN(x): Returns sign of value of x (-1 for negative, 0 for zero or 1 for positive)

POWER(x,y): Returns x to the power of y

SQRT(x): Returns square root of x

EXP(x): Returns exponent of x

LOG(y,x): Returns log of y base x

LN(x): Return natural log of x base e

LOG10(x): Returns log of x base 10

MOD(x,y): Returns remainder of x/y

RAND(): Returns a random number between 0 and 1

CEIL(x): Returns upper integer boundary of value x

FLOOR(x): Returns lower integer boundary of value x

ROUND(x,n): Rounds x value to n decimal places

TRUNC(x,n): Cuts off x value at n decimal places without rounding

String Functions

String functions are usually applied to attributes with string data types.

LENGTH(string): Returns number of characters of a string

UPPER(string): Returns upper case characters of the whole string

LOWER(string): Returns lower case characters of the whole string

TRIM(string): Removes trailing and leading whitespaces from string

TRIM(char FROM string): Removes trailing and leading specified char from string

LTRIM(string): Removes leading whitespaces from string

RTRIM(string): Removes trailing whitespaces from string

CONCAT(string1, string2, ...): Concatenate multiple strings into a single string

SUBSTRING(string, start, length): Returns a portion of string starting from start index for n length.

REPLACE(string, old_substring, new_substring): Replace old_substring from a given string with new_substring.

REGEXP 'Regular_Expression': Search for regular expression pattern

Finding index location of string syntax differs based on the software used.

INSTR(string, substring) : MySQL

CHARINDEX(substring, string, position) : SQL Server

POSITION(substring IN string) : PostgreSQL

Note: Regular expression syntax is best tested with the following website, instead of memorizing syntax: <https://regex101.com/>

Datetime Format Specifiers

%Y	4-digit year	%d	Day (1-31)
%y	2-digit year	%h	12 hours (1 – 12)
%m	Numeric month (1-12)	%H	24 hours (1 – 24)
%b	Abbreviated month (Jan – Dec)	%i	Minutes (0-59)
%M	Name of month (January – December)	%s	Seconds (0-59)

Note that datetime functions for different SQL engines have different syntax used. Thus, it is recommended to google it rather than memorizing syntax.

Null Functions

Null functions are triggered when there is a missing value from attributes and an alternative output is provided as part of the function.

COALESCE(attribute, alternative): Returns “alternative” if null value exist in attribute

Subqueries

Subqueries are queries nested inside another query, which can be found in various clauses like SELECT, FROM, WHERE or JOIN clause.

Subqueries must be surrounded by parenthesis () and including aliases except for subqueries in WHERE clause.

Subqueries are always executed from inner-most query to outer-most query.

Caveats to note for subqueries:

1. Using subqueries in SELECT clause requires result to return a single column and either zero or one row. If subquery returns one column and zero rows, the results will be filled with NULL values
2. Using subqueries in FROM clause requires the query to be independent without referring to values outside of the subquery.
3. Using subqueries in WHERE clause requires use of IN() or NOT IN() keywords if subquery returns multiple rows.

Advantages of subqueries in FROM clause:

1. Reduce problems to smaller steps required
2. Reduce table size in FROM clause to improve query performance

Advantages of subqueries in WHERE clause: Easier for code understanding

Correlated subqueries

Correlated subqueries are queries that refer to values obtained in outer queries.

Correlated subqueries execute once for every row from the outer query. Thus, correlated subqueries come with read performance issues.

It is recommended to use other approaches like Joins or CTE (Common table expressions) when possible.

SELECT clause

1. Select specified columns
SELECT c1, c2, ... FROM table1
2. Select all columns
SELECT * FROM table1
3. Select expressions
SELECT UPPER(c1), ROUND(c1,2) FROM table1
4. Select qualifying columns from multiple tables
**SELECT table1.c1, table2.c2 FROM table1
INNER JOIN table2 ON table1.c3 = table2.c3**
5. Select subqueries
**SELECT c1, (SELECT COUNT(*) FROM table2 WHERE table1.c2 =
table2.c2) AS count_name
FROM table1**
6. Select all distinct values from columns
SELECT DISTINCT c1, c2, ... FROM table1
7. Select number of distinct values from columns
SELECT COUNT(DISTINCT c1, c2, ...) FROM table1

Note that selecting number of distinct values from multiple columns is currently supported in MySQL and PostgreSQL.

FROM clause

1. Retrieve data from single table
SELECT * FROM table1
2. Retrieve data from multiple tables
**SELECT a.*, b.* FROM table1 a
INNER JOIN table2 b ON a.cl1 = b.cl1**
3. Retrieve data from table derived from subquery
**SELECT cl1, total
FROM
(SELECT cl1, cl2, COUNT(*) AS total FROM table2 GROUP BY cl1, cl2
HAVING COUNT(*)>5) A**

WHERE clause

1. Filters table based on single predicate (conditional statement)
SELECT * FROM table1 WHERE condition1
2. Filters table based on multiple predicates
SELECT * FROM table2 WHERE condition1 AND condition2
3. Filters table based on subquery
**SELECT cl1, cl2 FROM table1 WHERE cl3 IN (SELECT cl3 FROM table1
WHERE condition)**

GROUP BY clause

GROUP BY keyword used for summarizing rows into various groups, where collection of rows is specified in GROUP BY clause and aggregation functions are used in SELECT clause.

SELECT cl1, cl2, SUM(cl3)

FROM table1

GROUP BY cl1, cl2

Note that all non-aggregated columns listed in SELECT clause must be included in GROUP BY clause and at least one aggregation function must be included in SELECT clause for query to work.

HAVING clause

HAVING keyword used for filtering rows based on aggregation functions only after GROUP BY clause is used.

SELECT c11, c12, SUM(c13) as total

FROM table1

GROUP BY c11, c12

HAVING SUM(c13) >5

Note that filtering rows based on non-aggregated columns must be used in WHERE clause and using aliases on HAVING clause will result in error due to order of SQL statement execution.

ORDER BY clause

ORDER BY keyword used for specifying order of results from query to be sorted.

Default sorting is in ascending order (ASC), unless specify DESC for descending order.

SELECT c11, c12, ROUND(c14, 3) AS rounder

FROM table1

ORDER BY c13 ASC, c12 DESC, rounder ASC

Columns with aliases, other columns and expressions not within the SELECT clause can be used for sorting since FROM clause is executed first and SELECT clause is executed before ORDER BY clause.

Note that ORDER BY clause cannot be used within a subquery.

Limiting view of number of rows

Different syntax is used for this task depending on SQL software used.

For MySQL, PostgreSQL (Note that LIMIT clause is executed last after ORDER BY clause):

SELECT cl1, cl2

FROM owner

LIMIT n

For SQL Server:

SELECT TOP n cl1, cl2

FROM owner

For Oracle:

SELECT cl1, cl2

FROM owner

WHERE ROWNUM <= n

CASE statements

Case statements are similar to if-then-else statements.

CASE

WHEN condition1 THEN result1

WHEN condition2 THEN result2

ELSE result3

END AS Alias_Name

Functions for summarizing rows into lists

Using together with GROUP BY clause, the following functions can be used for summarizing multiple rows into list of values:

MYSQL	GROUP_CONCAT (DISTINCT x SEPARATOR ‘,’ ORDER BY x)
SQL Server	STRING_AGG(x, ‘,’) WITHIN GROUP (ORDER BY x)
PostgreSQL	ARRAY_AGG(DISTINCT x ORDER BY x)

Note the following caveats:

1. SQL server currently does not support “Unique” list
2. PostgreSQL only accepts ‘,’ as default separator
3. SQL server requires input of separator, while input of separator is optional for MYSQL

Complex Group-By Functions

Additional summary information can also be included in queries using the following complex group-by functions:

1. GROUP BY ROLLUP (a,b, ...,e)
 - Group by sets of a – e, a – d, a and ()
2. GROUP BY CUBE (a,b, ..., e)
 - Group by all combination of sets and ()
3. GROUP BY GROUPING SETS(a,b,...,())
 - Group by individual categories and (), which represents group function on entire table

Note that these functions are mostly supported in Oracle, PostgreSQL and SQL server.

MySQL only supports rollup function using the following syntax:

GROUP BY a, b, ..., e WITH ROLLUP

Window Functions

Window functions return value for each row of the dataset that relates to result from current query.

Window functions are mostly used in SELECT & ORDER BY clause and useful for combining non-aggregated columns and aggregated columns together without relying on GROUP BY clause.

Filtering columns related to window functions requires use of subqueries, because window functions are executed 2nd last prior to ORDER BY clause.

Basic syntax:

AGG_FUNC() OVER (PARTITION BY c1 ORDER BY c2)

Note that PARTITION BY clause has similar behavior as GROUP BY clause and basic aggregation functions can also be used in the syntax above.

By default, ORDER BY clause within OVER clause has “Cumulative function” behavior (hidden syntax is **RANGE BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW**)

The behavior (sliding windows) of ORDER BY clause within OVER clause can be adjusted to the following:

1. Moving N previous records
RANGE BETWEEN N PRECEDING AND CURRENT ROW
2. Moving N following records
RANGE BETWEEN CURRENT ROW AND N FOLLOWING
3. Moving N previous and N following records
RANGE BETWEEN N PRECEDING AND N FOLLOWING

Note that **RANGE BETWEEN** can be replaced by **ROW BETWEEN** instead.

RANGE BETWEEN will return same running value for each duplicated time, while ROW BETWEEN returns different running value for each row

List of window functions:

1. **ROW_NUMBER()** : Returns row number of every record without “ties”
2. **RANK()** : Returns rank of current row with “ties” allowed and “gaps” allowed
3. **DENSE_RANK()** : Returns rank of current row with “ties” allowed and “gaps” not allowed
4. **FIRST_VALUE(cl1)** : Returns value of cl1 at 1st row of window frame
5. **LAST_VALUE(cl1)** : Returns value of cl1 at last row of window frame
6. **NTH_VALUE(cl1, n)** : Returns value of cl1 at nth row of window frame
7. **LEAD(cl1, n)** : Returns value of cl1 at n rows after current row (Default n is 1)
8. **LAG(cl1, n)** : Returns value of cl1 at n rows before current row (Default n is 1)
9. **NTILE(n)** : Returns rank of current row by dividing partition into n-sections