

SQL & Python Summary

Rio Anggara Sufilin

What is SQL?

Structured Query Language (SQL) is a language used to communicate to database. It is the language that translates all that data into understandable information.

A **database** is a system that allows data to be easily stored, organized and managed.

Database Management System (DBMS)

Database Management System (DBMS) is a software for storing and retrieving users' data while considering appropriate security measures.

DBMS - Most Popular Database Management Systems

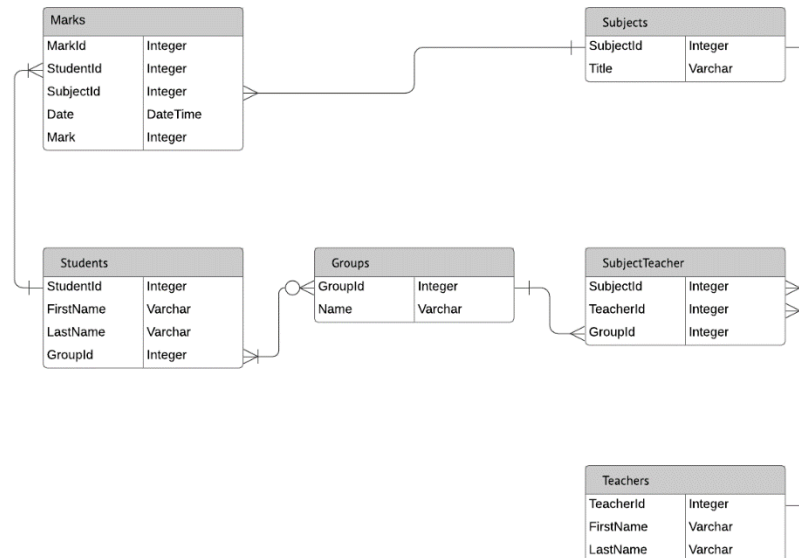


Relational Database

A **relational database** is a structural form of database that stores data in tables, and these tables can be somehow linked to each other. In the example of our company database, the employee table can be linked to a department table. The **relation** here is that the employee belongs to a department.

Database ER Diagram Example (Crow's Foot)

Lindi H | June 27, 2019



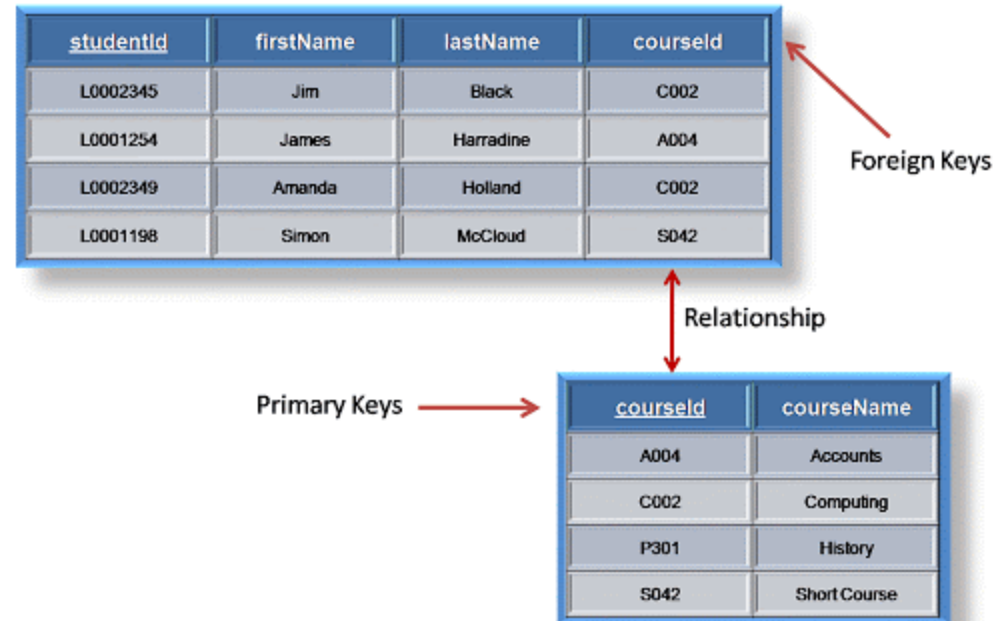
Database Terminologies

At its simplest, a database is made of **rows (records)** and **columns (field)**, similar to a spreadsheet, but it's far more powerful and has a vast number of features. Data is categorized and stored in the form of tables.

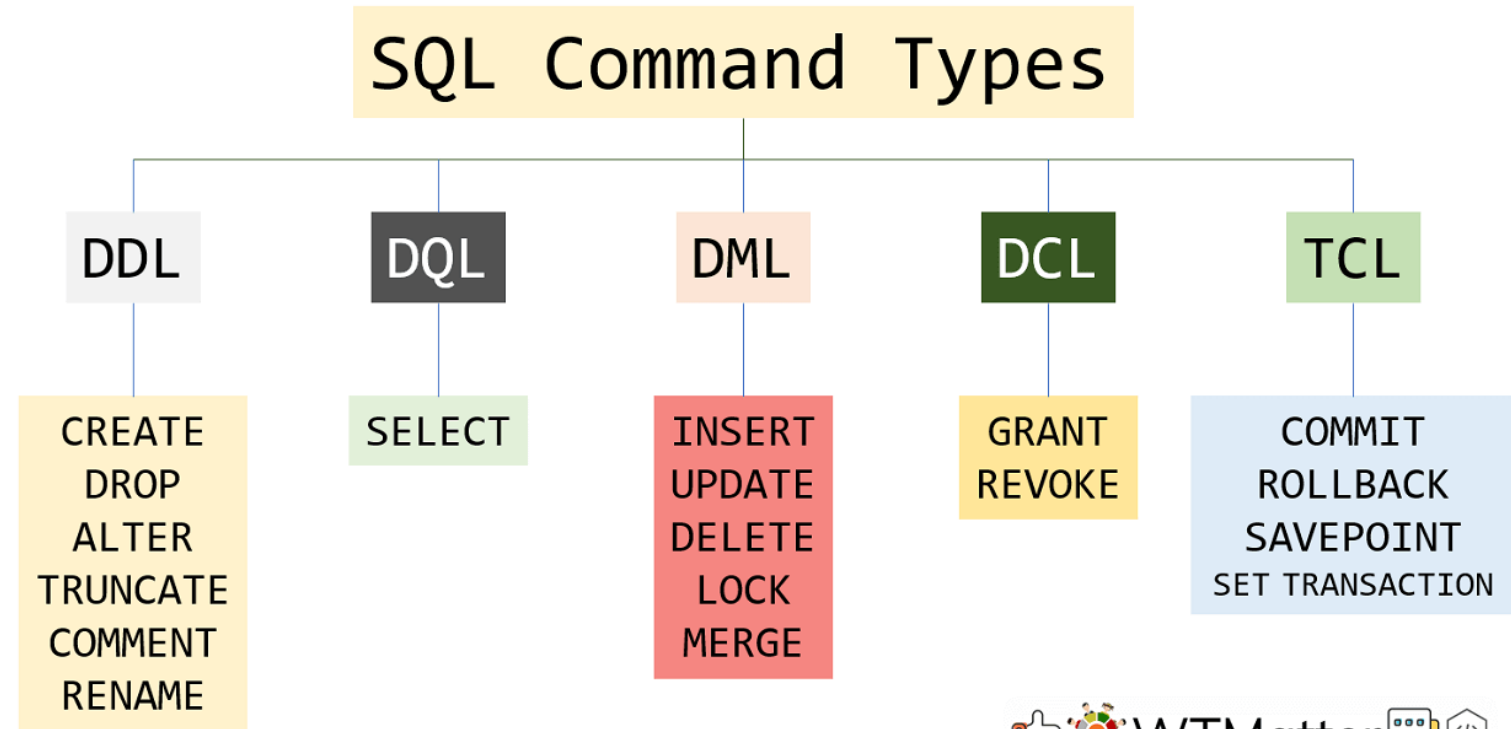
1	7	LINDA	WILLIAMS	LINDA.WILLIAMS@sakilacustomer.org	1
2	8	BARBARA	JONES	BARBARA.JONES@sakilacustomer.org	1
1	9	ELIZABETH	BROWN	ELIZABETH.BROWN@sakilacustomer.org	1
2	address_id	address	district	city_id	postal_code
1	5	1913 Hanoi Way	Nagasaki	463	35200
2	6	1121 Loja Avenue	California	449	17886
2	7	692 Joliet Street	Attika	38	83579
1	8	1566 Inegl Manor	Mandalay	349	53561
2	9	53 Idfu Parkway	Nantou	361	42399
1	10	1795 Santiago de Compostela Way	Texas	295	18743
2	11	900 Santiago de Compostela Parkway	Central Serbia	280	93896
2	12	478 Joliet Way	Hamilton	200	77948
1					1
2					0
1	21	DONNA	THOMPSON	DONNA.THOMPSON@sakilacustomer.org	1

Database Terminologies

- A **primary key** is used to ensure data in the specific column is unique.
- A **foreign key** is a column or group of columns in a relational database table that provides a link between data in two tables.

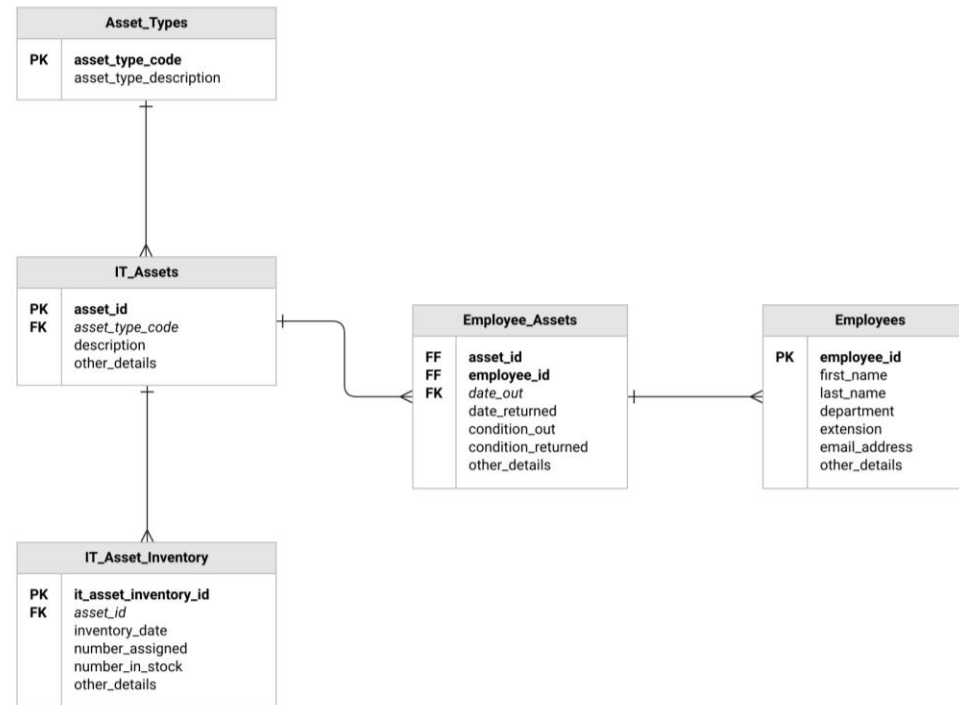


SQL Command Types



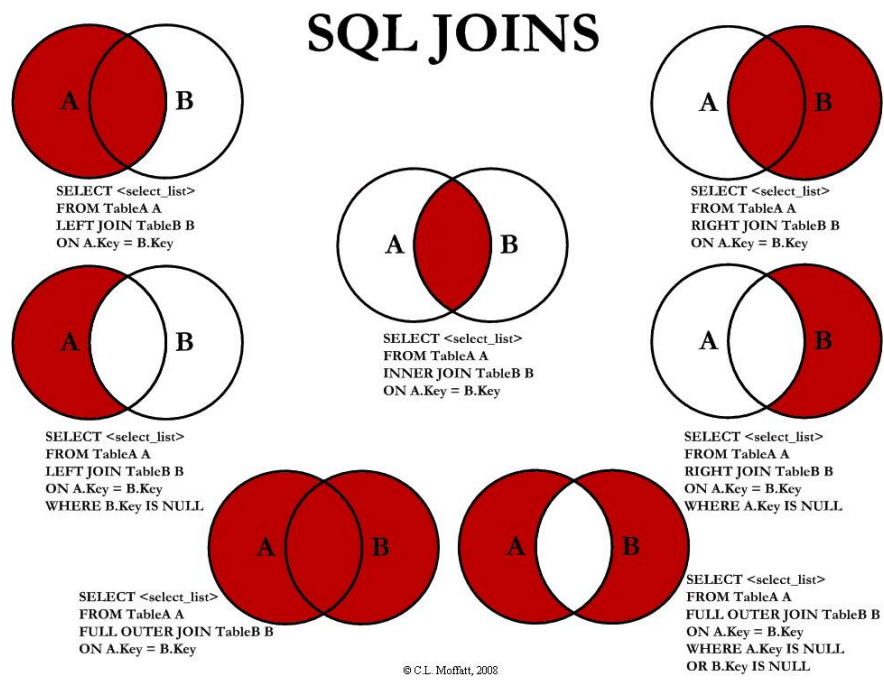
Entity Relationship Diagram

Entity Relationship Diagram, also known as ERD, ER Diagram or ER model, is a type of structural diagram for use in database design. An ERD contains different symbols and connectors that visualize two important information: **The major entities within the system scope**, and the **inter-relationships among these entities**.



SQL JOINS

JOINS in SQL are commands which are used to combine rows from two or more tables, based on a **related column between those tables**. There are predominantly used when a user is trying to extract data from tables which have **one-to-many** or **many-to-many relationships** between them



Python is a high-level, versatile, object-oriented programming language. It is useful and powerful while also being **readable** and **easy to learn**. This makes it suitable for programmers of all backgrounds and is likely the reason **Python** is **one of the most widely used programming languages (as of 2020)**.

What is
Python?

- Python works on different platforms (Windows, Mac, Linux, Raspberry Pi, etc).
- Python has a simple syntax similar to the English language.
- Python has syntax that allows developers to write programs with fewer lines than some other programming languages.
- Python runs on an interpreter system, meaning that code can be executed as soon as it is written. This means that prototyping can be very quick.
- Python can be treated in a procedural way, an object-oriented way or a functional way.

Why Python?

Python Installation

Many PCs and Macs will have python already installed.

To check if you have python installed on a Windows PC, search in the start bar for Python or run the following on the Command Line (cmd.exe):

```
C:\Users\Your Name>python --version
```

To check if you have python installed on a Linux or Mac, then on linux open the command line or on Mac open the Terminal and type:

```
python --version
```

Let's write our first Python file, called *helloworld.py*, which can be done in any text editor.

```
print("Hello, World!")
```

The way to run a python file is like this on the command line:

```
C:\Users\Your Name>python helloworld.py
```

The output:

```
Hello, World!
```

Python Hello World

To test a short amount of code in python sometimes it is quickest and easiest not to write the code in a file. This is made possible because Python can be run as a command line itself.

Type the following on the Windows, Mac or Linux command line:

```
C:\Users\Your Name>python
```

From there you can write any python, including our hello world example from earlier:

```
C:\Users\Your Name>python
Python 3.6.4 (v3.6.4:d48eceb, Dec 19 2017, 06:04:45) [MSC v.1900 32 bit (Intel)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> print("Hello, World!")
Hello, World!
```

Python Hello World

Python has no command for declaring a variable.

A variable is created the moment you first assign a value to it.

```
x = 5  
y = "John"  
print(x)  
print(y)
```

Python Variables

In programming, data type is an important concept.

Variables can store data of different types, and different types can do different things.

Python has the following data types built-in by default, in these categories:

Text Type:	<code>str</code>
Numeric Types:	<code>int</code> , <code>float</code> , <code>complex</code>
Sequence Types:	<code>list</code> , <code>tuple</code> , <code>range</code>
Mapping Type:	<code>dict</code>
Set Types:	<code>set</code> , <code>frozenset</code>
Boolean Type:	<code>bool</code>
Binary Types:	<code>bytes</code> , <code>bytearray</code> , <code>memoryview</code>

Python Data Types

Python Operators

Python divides the operators in the following groups:

- Arithmetic operators (+, -, *, /, etc.)
- Assignment operators (=, +=, -=, *=, etc.)
- Comparison operators (==, !=, >, <, etc.)
- Logical operators (*and*, *or*, *not*)
- Identity operators (*is*, *is not*)
- Membership operators (*in*, *not in*)
- Bitwise operators (&, |, ^, ~, etc.)

- The list is changeable, meaning that we can change, add, and remove items in a list after it has been created.

```
thislist = ["apple", "banana", "cherry"]  
print(thislist)
```

- Tuples are unchangeable, meaning that we cannot change, add or remove items after the tuple has been created.

```
thistuple = ("apple", "banana", "cherry")  
print(thistuple)
```

- Set items are unordered, unchangeable, and do not allow duplicate values. Sets are unchangeable, meaning that we cannot change the items after the set has been created, but you can add new items

```
thisset = {"apple", "banana", "cherry"}  
print(thisset)
```

- Dictionaries are used to store data values in key:value pairs. Dictionaries are written with curly brackets, and have keys and values

```
thisdict = {  
    "brand": "Ford",  
    "model": "Mustang",  
    "year": 1964  
}  
print(thisdict)
```

Python List, Tuple, Set, Dictionary

Python supports the usual logical conditions from mathematics:

- Equals: `a == b`
- Not Equals: `a != b`
- Less than: `a < b`
- Less than or equal to: `a <= b`
- Greater than: `a > b`
- Greater than or equal to: `a >= b`

These conditions can be used in several ways, most commonly in "if statements" and loops.

An "if statement" is written by using the if keyword.

```
a = 33
b = 200
if b > a:
    print("b is greater than a")
```

Python If ...
Else

Python has two primitive loop commands:

- **while** loops execute a set of statements as long as a condition is true.

```
i = 1
while i < 6:
    print(i)
    i += 1
```

- **for** loops is used for iterating over a sequence (that is either a list, a tuple, a dictionary, a set, or a string). This is less like the for keyword in other programming languages, and works more like an iterator method as found in other object-orientated programming languages. With the for loop we can execute a set of statements, once for each item in a list, tuple, set etc.

```
fruits = ["apple", "banana", "cherry"]
for x in fruits:
    print(x)
```

Python Loops

A function is a block of code which only runs when it is called.

You can pass data, known as parameters, into a function.

A function can return data as a result.

```
def my_function():  
    print("Hello from a function")
```

- A parameter is the variable listed inside the parentheses in the function definition.
- An argument is the value that is sent to the function when it is called.

```
def my_function(fname, lname):  
    print(fname + " " + lname)  
  
my_function("Emil", "Refsnes")
```

Python Functions

pandas is a software library written for the Python programming language for data manipulation and analysis.

pandas can be installed using:

```
pip install pandas
```

and can be used using :

```
import pandas as pd
```

Pandas is able to read several different types of stored data, including CSVs (comma separated values), TSVs (tab separated values), JSONs (JavaScript Object Notation, HTML (Hypertext Markup Language), among others.

Pandas

Data Wrangling with pandas Cheat Sheet

<http://pandas.pydata.org>

Syntax – Creating DataFrames

	a	b	c
1	4	5	6
2	7	8	9
3	10	11	12

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index=[1, 2, 3])
```

Specify values for each column.

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=["a", "b", "c"])
```

Specify values for each row.

	a	b	c
1	4	5	6
2	7	8	9
3	10	11	12

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d', 1), ('d', 2), ('e', 2)],
        names=["n", "v"])
```

Create DataFrame with a MultiIndex

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={
          "variable": "var",
          "value": "val"})
      .query("val >= 200"))
```

Tidy Data – A foundation for wrangling in pandas

In a tidy dataset:

Each variable is saved in its own column

Each observation is saved in its own row

Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

M * A

Reshaping Data – Change the layout of a data set

pd.melt(df)

Gather columns into rows.

df.pivot(columns="var", values="val")

Spread rows into columns.

pd.concat([df1, df2])

Append rows of DataFrames

pd.concat([df1, df2], axis=1)

Append columns of DataFrames

df.sort_values("mpg")

Order rows by values of a column (low to high).

df.sort_values("mpg", ascending=False)

Order rows by values of a column (high to low).

df.rename(columns={"year": "y"})

Rename the columns of a DataFrame

df.reset_index()

Reset the index of a DataFrame

df.reset_index()

Reset index of DataFrame to row numbers, moving index to columns.

df.drop(columns=["length", "height"])

Drop columns from DataFrame

Subset Observations (Rows)

df[df.length > 7]

Extract rows that meet logical criteria.

df.drop_duplicates()

Remove duplicate rows (only considers columns).

df.head(n)

Select first n rows.

df.tail(n)

Select last n rows.

df.sample(frac=0.5)

Randomly select fraction of rows.

df.sample(n=10)

Randomly select n rows.

df.iloc[10:20]

Select rows by position.

df.nlargest(n, "value")

Select and order top n entries.

df.nsmallest(n, "value")

Select and order bottom n entries.

Subset Variables (Columns)

df[["width", "length", "species"]]

Select multiple columns with specific names.

df["width"] or df.width

Select single column with specific name.

df.filter(regex="regex")

Select columns whose name matches regular expression regex.

regex (Regular Expressions) Examples	
"/./	Matches strings containing a period "."
"length"	Matches strings ending with word "length"
"^sepal"	Matches strings beginning with the word "sepal"
"^x[1-9]\$"	Matches strings beginning with 'x' and ending with 1, 2, 3, 4, 5
"^([species])\$"	Matches strings except the string "species"

df.loc[:, "x2": "x4"]

Select all columns between x2 and x4 (inclusive).

df.iloc[:, [1, 2, 5]]

Select columns in positions 1, 2 and 5 (first column is 0).

df.loc[df["a"] > 30, ["a", "c"]]

Select rows meeting logical condition, and only the specific columns.

	Logic in Python (and pandas)	
x	is less than	df.column.isin(values)
x	is equal to	pd.isnull(df)
x	is less than or equal to	pd.notnull(df)
x	is greater than or equal to	df.isnull().any()

http://pandas.pydata.org/ This cheatsheet inspired by R tidyverse data wrangling cheat sheet: <http://www.tidyverse.org/articles/2016/05/data-wrangling-cheat-sheet/> Written by: <http://www.tidyverse.org/>

Pandas Cheatsheet