

Data cleaning

④ There are 2 reason for using SQL to cleaning data.

(1) Data primarily stay in database, so in 1st level

(2) if we do something it only possible with SQL.

(2nd) SQL is really fast with huge data.

These are the 2 reason we have to know SQL.

⑤ In data wrangling there are 3 process:

1st) gathering data

2nd) Assessing data →

3rd) Cleaning data

⑥ Dat Assessing have 2 way to do?

① manually by scanning data with your eyes,

② through Programming

String Data type

From data types we can create 3 primary data types

① Numerical (int, decimal)

② Text (as string)

(char*)

③ temporal (Date and time)

String data type

① CHAR: The data type is used to store fixed-length string.

① if i make a char[10] (10 space will consume)

② if i store [Hello] in that char, it will consume 5 space,

since we specify [10] space, so it will consume 10 space

by adding space to [Hello -----] like this.

② length 255

③ we can't store more than 10 char, since we specify only up to 10 characters are allowed.

④ where we can use CHAR when we want to store fixed number of data type like {Phone number}

② VARCHAR:

① in varchar we specify a upper limit like 255

if we store 1 word it will occupy 1 space

it is actually dynamic

② varchar generally use for those column, we don't the lengths of the data. like: Name, email,

③ TEXT:

① This data type is used to store larger amounts of variable-length string data than Varchar.

② it can store up to 65,535 character.

③ used in text data, such as blog post or comments.

④

medium-text :

- ① This data type is used to store even larger amount of text data than text.
- ② It can store up to 16,777,215 characters.
- ③ Generally used to in {long-form articles on legal documents} Privacy-Policy, terms and conditions}

⑤

long-text :

- ① This data type is used to store the largest amount of text data.
- ② It can store up to 4,294,967,295 characters.
- ③ Used in entire books or large collection of data.

wildcards

- ① The Like operator in MySQL is used to match a string value against a pattern using wildcard characters. It is commonly used in select, where and join clause to filter or join rows based on a pattern match.

Like operator

wildcards

Percentage (%)

% thin represent zero, one or more characters,

a single character

white (-)

underscore (-)

④ find all the movie name, that has only 5 character.

Select * from movies

where name Like '_____';

{ one underscore represent
one word/char }

⑤ find all the movie that start with A and has 5 character.

Select * from movies

where name Like 'A_____';

⑥ find all the movie name that has (man) word?

Select * from movies;

where name like '% man %'

⑦ find all the movie name that start with (man)word?

Select * from movies;

where name like '(man %)'

④ Find all the movies that end with 'man' word;

Select * from movies

where name like '%man';

String function

→ Upper/lower operator

Select name, upper(name), lower(name) from movies;

concat_ws (with separator)

→ Concat and

⑤ I want to write movie name and director name?

Select concat(name, ' => ', director) from movies;

⑥ also adding star name

Select concat(name, ' => ', director, ' => ? ', star) from movies.

- ④ In SQL there are 2 types of type conversion are there
- ↳ Implicit (Internally understand (0) this number need to convert into string in order to concate)
 - ↳ Explicit (using cast function)

Concat_ws (with separator)

① Select concat_ws(' => ', name, director, star) from movies;

Substring → last 5 characters

② I need to extract 5 characters from movie name?

Select name, substr(name, 1, 5) from movies;

③ If we don't provide end value, it will automatically select last index of string.

Select name, substr(name, 4) from movies;

⑩ we can even put negative value.

⑪ so if i want first character from last five character

Select name, substr(name, -5, 1) from movies;

⑫ replace string operator

① Select Replace("hello world", "world", "India")

my string ↓
the substring
I want to replace

② change the movie name that has a substring, man change it with women.

Select Replace(name, 'man', 'woman') from movies.

⑬ Reverse function

⑭ find those movie that has a Palindrome word

Select name from movies

where name = reverse(name)

⑪ Char-length vs length

both are not same, the main difference between CHAR-length and length is that CHAR-length return the length of string. whereas length return the length of a string in bytes.

② both are same, return character length, but there is a subtle difference.

Café → if apply length then i will get 5
→ but if we apply char-length " " 4

⑪ insert (str, position, length, newstr)

str: the original string to insert into.

position: the position at which to insert the new substring. The first position is 1.

length: the number of character to replace

'newstr': the new substring to insert.

Select insert ("Hello world", 7, 0, "india") → Hello indiaWorld
select insert ("Hello world", 7, 5, "india") → Hello indice

Left Right same as substring

→ the shining ←

Select name, Left(name, 3), Right(name, 3) from movies.

• Left \Rightarrow the. (first word there \rightarrow 3 word)

Right \Rightarrow ~~sh~~ iny (last word there \leftarrow 3 word)

Repeat operation

Select Repeat(name, 3) from movies;

Trim [trim and rtrim]

Trim Remove Space from first and last.

Select RTrim("nitish")

\Rightarrow nitish.

④ we can also remove any character ?

Select Trim(Both (.) from ".....nitish-----")

\Rightarrow nitish.

④ if want ~~all~~ character Shanto

Select Trim(Leading " ") from "... Shanto")

⑤ if i want ~~last~~ character Shanto -- --

Select Trim(Trailing ".") from (" Shanto -- --")

⑥ LTRIM and RTrim only Remove Space.

LTRIM(" Shanto -- -- "); → Shanto -- --

RTRIM(" Shanto -- -- "); → -- -- Shanto.

substring_index (Python split function) super important

- ① Select `substring_index("www.campusx.in", ".")`
⇒ www

first occurrence (2nd)
string result from
- ② Select `substring_index("www.campusx.in", ".", 2)`
⇒ www. campusx.
- ③ `substring_index("www.campusx.in", ".") =`
strcmp (string compare)
→ in.

the strcmp() function return an integer that indicates the relationship between the two strings.

 - ④ if str1 is less than str2, the function return a negative integer
 - ⑤ if str1 is greater than str2, the function return a positive int.
 - ⑥ if same str1 = str2 then 0

select strcmp("mumbai", "Delhi") \Rightarrow 1

(0, "mumbai") < (0, "Delhi")
m comes first in ASCII.

② locate ("Hello, world")

\Rightarrow locate actual field where the index of word starts.

select locate("w", "hello, world") \Rightarrow 7.

③ LPAD and RPAD

doing padding. suppose this is phone number
if want add +880,

1811146200
10

select LPAD("1811146200", 14, "+880")

select RPAD("1811146200", 14, "+880")