

Date functions in SQL :-

NOW() :- Returns current date and time

CURDATE() :- current date

CURTIME() :- current time

DATE() :- extract the date part of a date(s) date/time expression

EXTRACT() :- returns single part of a date(s) date/time

DATE ADD() :- adds a specified ~~part~~ time interval to a date

DATE SUB() :- subtracts a time interval from a date

DATE DIFF() :- returns the number of days b/w two dates

DATE FORMAT() :- displays date/time data in different formats

TO_DATE() :- converts a character string to a date

Year :- extracts the year component from a date
Month :- month

String Functions :-

CONCAT() :- (str1, str2, ...) concatenates

SUBSTR(str, start, length) :- extracts a substring from a string starting at its specified position and with the specified length

SPLIT() :- (str, delimiter) :- splits a string into an array of substrings based on a delimiter.

REPLACE(str, search, replace) :- replaces all occurrences of a specified substring with another substring in a string.

REGEXP_REPLACE(str, pattern, replace) :- replaces all occurrences of a specified regular expression pattern with a replacement string in a string.

INSTR(str, substr) :- returns the position of the first occurrence of a substring within a string.

CONCAT_WS(delimiter, str1, str2, ...) :- concatenates two or more strings together with a specified delimiter b/w them

INITCAP(str) :- convert the first character of each word in a string to uppercase and the lower case

TRIM() :- removes leading and trailing white space from a string

REVERSE(str) :- reverse the characters in a string

TRANSLATE(str, from, to) :- Replaces character mapping based on from, to

on a mapping defined by two strings.

• FIND-IN-SET (str, strlist) :- Returns the position of string with a comma-separated list of strings.

different function in ~~Python~~ Python:-

print() :- outputs a specified message or value to the screen (*) :- returning the depth (number) of objects such as a string, list, or tuple.

type() :- datatype of an object

range() :- generates a sequence of numbers with a specified range.

input() :- prompts the user to enter input from the console.

int() :- converts a value to an integer data type.

str() :- Converts a string to a string.

list() :- iterable object into a list

max() :- returning from largest items from a collection of data.

min() :- returning smallest item.

sum() :- returning the sum of all items in an iterable.

abs() :- returning absolute value of a number.

round() :- Rounds a number to specified number decimal places.

zip() :- combines multiple iterables into a single iterable of tuples.

map() :- Applies a given function to each items in an iterable and returns an iterator of the outputs.

filter() :- Returns an iterator containing only the items from an iterable that satisfy a given condition.

open() :- opens a file for reading, writing (both).

dir() :- return list of names in current namespace.

enumerate :- return iterator that generates tuple, containing index and value of each item in an iterable.

different function of PANDAS API :-

1. Data loading and input :-

`read_csv()` :- Read csv (comma-separated) files into a Data frame.

`read_excel()` :- Read Excel files into a Data frame.

`read_sql()` :- SQL query or database table into a Data frame.

`read_json()` :- Read json files into a Data frame (JavaScript object notation)

2. Data exploration & manipulation :-

`head()` :- Returns the first n rows of a Dataframe

`tail()` :- Returns last n rows

`info()` :- provide summary of a Dataframe's structure and data types

`describe()` :- Generates descriptive statistics of DF

`columns()` :- Returns column names labels of DF

`unique()` :- Unique values in a column

`sort_values()` :- Sorts a DF by one or more columns

3. Data Selection and filtering :-

`loc[]` :- Accesses rows & columns by label

`iloc[]` :- integer based indexing

`isin()` :- filters row based on whether a column value is in given list

`query()` :- filters row using a Boolean expression.

4. Data Aggregation and grouping :-

`groupby()` :- Groups data based on one or more columns

`agg()` :- Applies aggregation functions to grouped data

`pivot_table()` :- Create spreadsheet-style pivot table based on grouped data.

5. Data Cleaning and Transformation :-
- `fillna()` :- fills missing values in a Data frame with a specified method.
 - `drop_duplicates()` :- removes duplicate row from Data frame.
 - `replace()` :- replaces values in DF with specified values.
 - `apply()` :- applies a function to each element of row or column of a Data frame.

6. Data Visualization :-
- `plot()` :- generates various types of plots, such as line plots, bar plots, etc.
 - `hist()` :- creates a histogram from a Data frame or series.
 - `boxplot()` :- creates a box plot from a D.F or series.

DF is the NUMPY API :-

1. Array creation :-

`numpy.array()` :- creates an array from a python list or tuple.

`numpy.zeros()` :- an array filled with zeros.

`numpy.ones()` :- creates an array filled with ones.

`numpy.empty()` :- creates an uninitialized array.

`numpy.arange()` :- creates an array with evenly spaced values.

`numpy.linspace()` :- creates an array with evenly spaced values in a specified range.

`numpy.random.rand()` :- creates an array of random numbers b/w 0 & 1.

2. Array manipulation :-

`numpy.reshape()` :- Reshape an array into a specified shape.

`numpy.concatenate()` :- Concatenates arrays along a specified axis.

`numpy.split()` :- Splits an array into multiple sub-arrays

`numpy.transpose()` :- Transposes an array.

3. Mathematical functions :-

`numpy.sum()` :- Computes the sum of array elements

`numpy.mean()` :- Computes the arithmetic mean of array elements

`numpy.min()` :- finds the minimum value in an array

`numpy.max()` :- " " " max " "

`numpy.sin()`, `numpy.cos()`, `numpy.tan()`, compute trigonometric functions.

`numpy.exp()` :- computes the exponential of all elements in array.

`numpy.log()` :- computes the natural logarithm of all elements in an array.

4. Array operations :-

`numpy.dot()` :- computes the dot product of 2 arrays

`numpy.transpose()` :- Transposes an array.

`numpy.sort()` :- sorts the elements of an array.

`numpy.unique()` :- finds the unique elements in an array

`numpy.argmax()` :- Returns the indices of the maximum value along a specified axis

5. Random Number Generations :-

`numpy.random.randint()` :- Generates random integers

`numpy.random.randn()` :- Generates an array of random numbers from the standard normal distribution,

`numpy.random.choice()` :- Generates random samples from a given 1-D array.

`numpy.random.shuffle()`: shuffles the elements of an array randomly.

* function of Spark :-

1. Data loading and input :-

`spark.read.csv()`: Reads CSV file into a DataFrame
`spark.read.parquet()`: Reads Parquet file into DataFrame
`spark.read.json()`: " JSON " " "
`spark.read.text()`: Read text files into DF
`spark.read.jdbc()`: Reads data from JDBC data source into DataFrame.

2. DataFrame operations :-

`Dataframe.select()`: selects specific columns from a DF
`Dataframe.filter()`: filters rows based on a condition
`Dataframe.withColumn()`: Adds or replaces a column in DF
`Dataframe.groupby()`: Groups data based on one or more columns

`Dataframe.join()`: Joins two Dataframes based on a column expression.

`Dataframe.sort()`: sorts the Dataframe by one or more columns.

3. Aggregation and window functions :-

`Dataframe.agg()`: applies aggregation functions to grouped data.

`Dataframe.groupby().agg()`: Aggregates data based on grouping.

`PYSPARK.SQL.functions.sum()`: Computes the sum of a column.

`PYSPARK.SQL.functions.count()`: counts the number of rows of non null values in a column

`PYSPARK.SQL.functions.avg()`: computes the average of a column.

`PYSPARK.SQL.function.rank()`: Computes the rank of a row within a partition.

4. Data Transformation and cleaning :-
Data frame, with column Renamed () :- Renames a column in a
Data frame.

Data frame, drop () :- Drops specified columns from a data frame

Data frame, fillna () :- fills missing values in Data frame
with specified value.

Dataframe, na, drop () :- Drops rows with missing values
from a Data frame.

PySpark, SQL functions, when :- Applies conditional transformation
to column.

5. SQL Queries :-

spark, sql () :- executes SQL queries on registered table
of Data frames.

Data frame, create or replace Temp View :- Register in Data frames
as temporary table.

spark, catalog, list tables () :- lists the table available
in the spark catalog.

spark, catalog, refresh table () :- Refreshes the metadata
of a table in spark catalog.

spark, catalog, refresh table () :- Refreshes the metadata of
a table in the spark catalog.

6. Output and storage :-

Data frame, write, CSV () :- writes a Dataframe to CSV file.

Data frame, write, Parquet () :- writes a Data frame to
Parquet file.

Data frame, JDBC () :- writes a Data frame to a JDBC Datasource

Data frame, write, JSON () :- " " " " " JSON file

~~It~~ functions in the RDD API is

Transformation functions

`map()`: Applies a function to each element and returns a new RDD

filter () :- filter the elements based on a condition

`flatmap()` - applies a function to each element and returns a new RDD by flattening the result

`union()` :- Returns The union of two RDDs.

distinct): Removes duplicates and returns a new RDD.

`SORT BY ()`: sort the ROD elements based on key

`groupByKey()` :- Groups the elements by key.

Reduce By key() :- Aggregation The Values of each key

`join()` - Joins two RDDs based on a common key
`coalesce()` - Reduces the number of partitions

Action function :-

`collect()` :- Returns all the elements of the RDD as an array.

`count()` :- returns the number of elements ⁱⁿ of the RDD.

`first()`, returning the first elements of the RDD.

`take()` :- Returns The first n elements of The RDD

`reduce()` :- Reduces the elements of the RDD using a binary operator for each:-

Appendix in the back

Applies a function to each element of the RDD pair RDD functions:-

key() :- Returns an RDD of the keys

values :- ' ', values

`sortByKey()` :- sorts the RDD elements by key.

Reduce by key() :- Aggregates The Values of

Each key ~~miss~~

`groupByKey()` :- groups The elements by Key.

persistence and caching functions:

cache() : caches the RDD in memory for faster access

persist() persists the RDD in memory & disk storage.

unpersist() : removes the RDD from memory (8) disk storage.