

Date functions in SQL

- 1) `GetDate()`: Returns current date and time.
- 2) `current_timestamp`: Returns current date and time.
- 3) `Datepart(part,date)`: extracts specific part of date or ^{time value}.
- 4) `Year(date)`: Returns year from date value.
- 5) `Month(date)`: Returns month from a date value.
- 6) `Day(date)`: Returns day from a date value.
- 7) `Dateadd(interval, num, date)`: Returns day from a date value.
Add specific no. of intervals to a date.
- 8) `Date diff(interval, start date, end date)`: Calculate Add a specific diff between two dates in specified interval.
- 9) `Convert(data-type, expression, style)`: Converts a date value to a specified data type.
- 10) `Format(date, format)`: formats a date value using a specified format.

String fun:

- 1) `Len(str)`: length of str.
- 2) `Upper(str)`: upper
- 3) `Lower(str)`: lower
- 4) `Left(str,len)`: Return specified number of characters from left side of string.
- 5) `Right(str,len)`: Right side of string
- 6) `Substr(str,start,len)`: Returns a substr from str starting at a specified position with a specified length.

7) `Replace(str, old, new)` : Replaces all occurrences of substr with in str with new substr.

8) `Concat(str1, str2)` : Concatenates

9) `Trim(str)` : Removes leading and trailing spaces from str.

10) `Charindex(substr, str)` : Returns the starting position of a sub string within a string.

Pandas API func

1) `Pd.read_csv()` : Reads data from csv file and creates dataframe

2) `df.head()` : Returns the first n rows of a dataframe (default 5)

3) `df.tail()` : Returns last n rows of dataframe

4) `df.shape` : Returns the dimensions (row, col) of a dataframe.

5) `df.info()` : Provides a concise summary of dataframe.

6) `df.describe()` : Generates descriptive statistics of a dataframe include count, min, max etc

7) `df.columns` : Returns the column labels of a dataframe.

8) `df.loc[]` : Access a group of rows and columns by label(s) or a boolean array.

9) `df.iloc[]` : Access a group of rows and integer position

10) `df.groupby()` : Groups data in a dataframe based on specific column

11) `df.merge()` : Combines 2 dataframes based on common columns

12) `df.sort_values()` : Sorts dataframe by one or more columns.

13) `df.dropna()` : Drop rows or columns with missing values from dataframe.

14) `df.fillna()` : Fills missing values in a dataframe with specified value.

15) `df.to_csv()` : converts a dataframe to a csv file.

Numpy API functions

- 1) np.array(): creates a numpy array from a python list or tuple
- 2) np.zeros(): creates an array filled with zeros.
- 3) np.ones(): creates an array filled with ones.
- 4) np.arange(): "an array with specified numbers evenly spaced values within a specified range."
- 5) np.linspace(): creates an array with specified numbers evenly spaced values within a specified range.
- 6) np.random.rand(): generates an array of random numbers between 0 and 1.
- 7) np.random.randint(): generates an array of random integers within a specified range.
- 8) np.reshape(): reshapes an array to a specified range.
- 9) np.reshape(): reshapes an array to a specified range.
- 10) np.mean(): computes mean of array.
- 11) np.sum(): sum of elements.
- 12) np.max(): max value in an array.
- 13) np.min(): min value in an array.
- 14) np.dot(): computes dot product of two arrays.
- 15) np.transpose(): transposes an array.

RDD API functions

1) Creation function:-

- * `SC.parallelize()`: creates an RDD from a local collection.
- * `SC.textFile()`: Reads a text file and creates an RDD where each line represents an element.

2) Transformation function:-

- * `map()`: Applies a transformation function to each element of RDD and returns a new RDD.
- * `filter()`: filters the elements of RDD based on a given Predicate function.
- * `flatMap()`: Applies a transformation function that returns an iterator for each element and flattens the results.
- * `union()`: Combines two RDD's in a single RDD.
- * `distinct()`: Returns a new RDD with distinct elements.
- * `sortByKey()`: Sorts the RDD elements based on a given key.
- * `groupByKey()`: Groups the values of each key in the RDD.
- * `reduceByKey()`: Reduces the values of each key using a specified reduction function.
- * `join()`: Performs an inner join between two RDD's based on a common key.

3) Action function:-

- * `count()`: Returns the number of elements in RDD.
- * `collect()`: Returns all elements of RDD as an array to client program.
- * `first()`: Returns the first element of RDD.
- * `take(n)`: Returns the first n elements of RDD.

- * reducer(): Applies a reduction function to aggregate the elements of RDD.
- * foreach(): Applies a function to each element of RDD by persistence function:-
- * cache(): persists the RDD in memory.
- * persist(): persists the RDD with a specified storage level.
- * unpersist(): removes the RDD from memory.
- * partitioning function:-
- * repartition(): Repartitions the RDD into a specified number of partitions.
- * coalesce(): Reduces the number of partitions in RDD.
- * other functions:-
- * sample(): Returns a random sample of RDD.
- * foreach partition(): Applies a function to each partition of RDD.
- * glom(): Returns RDD where each element is an array of elements from the original RDD partition.

Spark Session Object:-

It is an entry point and main interface for interacting with Apache Spark. It provides unified API for creating and working with Spark functionality.

Important methods and properties of SparkSession object:-

Creating a SparkSession:-

- * `SparkSessionBuilder`: creates a `SparkSessionBuilder` instance to configure and build a `SparkSession`.
- * `appName(name: String)`: Sets the name of Spark application.
- * `master(master: String)`: sets the master URL for connecting to the Spark cluster.
- * `Config(key: String, value: String)`: sets a configuration property for Spark.
- * `getOrCreate()`: Returns an existing `SparkSession` instance or create a new one if it doesn't exist.

Accessing properties:-

- * `SparkContent`: Returns the underlying `SparkContent`.
- * `Catalog`: provides access to catalog metadata, such as tables and databases.
- * `Conf`: Returns the `SparkConf` object, which holds Spark configuration properties.
- * `SessionState`: Returns the `SessionState` object, which manages the state of `SparkSession`.

Working with Dataframes and Datalets

- * `CreateDataFrame(data: RDD[Row], schema: StructType)`: Create a DataFrame from an RDD of Row objects and Schema.
- * `CreateDataFrame(data: java.util.List[Object], beanClass: Class[_])`: Creates a DataFrame from a Java list of objects and corresponding Java bean class.
- * `read()`: Returns a 'DataFrameReader' to read data from various sources.
- * `sql(sqlText: String)`: Executes a SQL query and returns the result as a DataFrame.
- * `table(tableName: String)`: Returns a DataFrame representing a table registered in the Catalog.

Executing Actions:

- * 'DataFrame' and 'Dataset' objects provide methods like `show()`, `Count()`, `Collect()` etc. to perform actions on data.

Managing Resources:

- * `stop()`: stops the 'Spark Session' and releases the allocated resources.
- * `close()`: Closes the Spark session.