

Functions of Hive

String:

length() - returns the length of string

reverse() - reverse input string

concat() - returns string or bytes after concatenating all strings

space() - returns a string with spaces

split & →

upper() | lower() - converts to upper case/
lower case

Date: current_date() - Current date is returned

current_timestamp() - Current system time
date in timestamp format

date_diff() - diff b/w specified dates

Year, Quarter, Month, Day

last_day()

functions of list:

append() - adds element at end of list

pop() - removes at specified position

sort() - sorts the list

reverse() - reverses the list

remove() - first item with specified value
is removed

insert() - adds element at specified position

index() - returns index of first element

extend() - add elements of list to end of

count() - returns no. of elements with value

copy() - returns a copy of list

clear() - removes all elements from list

Namenode: Has all the information such as which file is saved in which location, access-time of file, which user access a file on current time.

Secondary: Stores data when namenode fails it is used to restart the namenode.
It runs on different machines for memory management.

Functions in Python Pandas

- 1) `read_csv()` — helps read a comma separated value file into pandas dataframe. Path of file to be read is mentioned. Delimiters other than comma, like or tab.
- 2) `head()` — Used to return first n rows of a dataset.
`df.head()` — returns first 5 rows of dataframe
- 3) `describe()` — Used to generate descriptive statistics of data in pandas df
Quick overview of dataset
- 4) `memory_usage()` — returns a pandas series having memory usage of each column in a pandas

- 5) `astype()` - cast a Python object to a particular data-type.
- 6) `loc[:, :]` - helps to access a group of rows and columns in a dataset
- 7) `to_datetime()` - Convert Python object to datetime format
- 8) `value_counts()` - returns Pandas Series containing counts of unique values.
- 9) `drop_duplicates()` - returns Pandas dataframe removing duplicate rows
not a number
- 10) `fillna()` - helps to replace all NaN values in a dataframe by imputing these missing values with more appropriate values.

Numpy functions:

- 1) `np.array()` - create an array in Numpy
- 2) `np.mean()` - find mean value of array
- 3) `np.dot()` - find dot product of two arrays
- 4) `np.random.randint()` - create array with random integer values bw specified range

- 5) np.random.rand() - create an array with random values b/w 0 and 1
- 6) np.linspace() - create an array with a specified number of evenly spaced values.
- 7) np.ones() - create array filled with ones
- 8) np.zeros() - man value in array
- 9) np.zeros() - create array filled with zeros
- 10) np.median() - median value of array

Spark Session object

- 1) sparksession.builder() - TO create a new spark session
- 2) stop() - stop/end spark session
- 3) range() - returns single column with long type
- 4) Conf - returns runtime config object
- 5) Create data frame - data frame from collections
RDD

- 6) `CreateDataset()` — creates dataset from collection,
Dataframe & RDD
- 7) `emptyDataFrame()` — creates empty data frame
- 8) `catalog` — returns catalog object to access metadata
- 9) `get defaultSession()` — returns default spark session returned by builder
- 10) `getOrCreate()` — creates or returns a spark context

Functions in RDD API's

1) Transformation functions:

`map()` — Applies transformation function to each element of RDD & returns a new RDD

`filter()` — filters elements based on a given condition and return a new RDD containing filtered elements

`flatmap()` — applies a transformation function that returns an iterator for each element & flattens results into new RDD

`distinct()` — returns a new RDD containing distinct elements from original RDD

`sortBy()` — sorts elements of RDD based on specified criteria & returns a new RDD

2) Action functions:

collect() - returns all elements of RDD as an array to driver program

count() - returns no. of elements in RDD

reduce() - Aggregates elements of RDD using a specified function.

take() - Returns the first N elements from RDD as an array

foreach() - applies a function to each element of RDD.

3) Pair RDD functions:

reduceByKey() - Performs a reduction operation on values of a RDD

groupByKey() - groups values of a pair RDD based on key based on key

sortByKey() - sorts elements

join() - Performs an inner join b/w two pair RDD's based on

cogroup() - groups the values of multiple pair RDD's sharing same key

4) Persistence functions:

`cache()` - Persists RDD in memory for faster future access

`Persist()` - Allows specifying different storage levels for persisting RDD

`unpersist()` - Removes RDD from memory/disk storage

5) Input/Output functions:

`textfile()` - Reads a textfile and converts it into an RDD of strings

`savesTextfile()` - writes the contents of RDD to text file

Functions in Spark API's

i) Data Loading & I/O:

`Spark.read.csv()` - read CSV files into DF

`spark.read.parquet()` - read parquet files into DF

`spark.read.json()` - read JSON files into DF

`spark.read.text()` - read text files into DF

`Spark.read.jdbc()` - read data from JDBC data source to use into DF

2) Dataframe operations:

Dataframe.select() - selects specific columns from DF

Dataframe.filter() - filters row based on a condition

Dataframe.withColumn() - Adds or replaces a column in DF

Dataframe.groupBy() - groups data based one or more columns

3) Aggregation & window functions:

Dataframe.agg() - applies aggregation function to grouped data

Dataframe.groupby().agg() - Aggregates data based on grouping

Pyspark.sql.functions.sum() - computes sum of a column

Pyspark.sql.functions.count() - Counts no. of rows or non-null values in column

Pyspark.sql.functions.avg() - Computes average of column

Pyspark.sql.functions.rank() - Computes rank of a row within a partition

a) Data Transformation & Cleaning:

Dataframe. withColumnRenamed() - renames column in Df

Dataframe. drop() - drops specified column from Df

Dataframe. fillna() - fills missing values in a Df with a specified value

Dataframe. na. drop() - drops rows with missing value from Df

Dataframe. sql. functions. when() - applies conditions + transformations to a column

5) SQL Queries:

Dataframe. withColumnRenamed(F - Renames column in Df)

Spark. sql() - executes SQL queries on registered tables or Df

Dataframe. createOrReplaceTempView() - Registers a Df as temporary table

spark. catalog. listTables() - lists the tables available in spark catalog

Spark. catalog. refreshTable() - Refreshes metadata of a table in spark catalog