

1) np.array(): This fn is used to create a Numpy array from a python list or tuple. Converts i/p data into an ndarray object.

2) np.arange(): Returns an array with evenly spaced values within a specified range.

3) np.reshape(): Used to change the shape of the array without changing its data.

4) np.mean(): np.mean() computes arithmetic mean along a specified axis or the whole array. Calculates avg of elements in the array.

5) np.max(): Returns max value from an array along specified axis or whole array.

- 1) df.head() : returns first n rows of a dataframe
- 2) df.tail() : returns last n rows of a dataframe
- 3) df.describe() : generates descriptive stats of numerical columns in a dataframe
- 4) df.sort\_values() : used to sort dataframes by one or more columns.
- 5) df.groupby() : used to group rows of a dataframe based on one or more columns.
- i) df.drop() : remove rows/columns axis.

- 1) getdate(): returns current system date & time.
- 2) dateadd(): adds or subtracts a specified interval
- 3) datediff(): calculates diff between two dates.
- 4) datepart(): returns specified part of a date, such as date, month & year.
- 5) convert(): converts a given value to a specified data type such as converting date to a different format.
- 6) unix\_timestamp(): returns the unix epoch seconds from a given date.

## RDD API

↳ Represents immutable distributed collection of objects that can be processed in parallel across clusters.

↳ RDD API provides a set of ops.

i) map - applies a fn to <sup>Transform</sup> each element of RDD & returns a new RDD of results.

ii) filter. - operation used to selectively retain elements based on a condition.

iii) reduceBykey - performs reduce operation based on a key value set of values.

Performs inner join ← iv) join - combines 2 RDDs based on common key.

v) collect - retrieves all data from RDD as in-memory DS.

vi) count -

vii) saveAsTextFile. - saves the result <sup>RDD</sup> in a text file