

PJ LOKESH

211720104084

HOME WORK 9

1.##What is Apache Spark Streaming

Apache Spark Streaming is a scalable fault-tolerant streaming processing system that natively supports both batch and streaming workloads.

2.##Describe how Spark Streaming processes data?

Spark Streaming receives live input data streams and divides the data into batches, which are then processed by the Spark engine to generate the final stream of results in batches.

Spark Streaming provides a high-level abstraction called a discretized stream or DStream, which represents a continuous stream of data.

3.##What are DStreams?

Discretized Streams (DStreams):

A Discretized Stream or DStream is the basic abstraction provided by Spark Streaming. It represents a continuous stream of data,

Either the input data stream received from source, or the processed data stream generated by transforming the input stream.

4.##What is a StreamingContext object?

Public class StreamingContext extends Object implements Logging. Main entry point for Spark Streaming functionality.

It provides methods used to create DStreams from various input sources. It can be either created by providing a Spark master URL and an app name, or from `org.apache`.

5.##What are some of the common transformations on DStreams supported by Spark Streaming?

It applies one every batch meaning every RDD in a DStream. It includes common RDD transformations like `map()`, `filter()`, `reduceByKey()` etc.

Although these functions seem like applying to the whole stream, each DStream is a collection of many RDDs (batches). As a result, each stateless transformation applies to each RDD.

6.##What are the output operations that can be performed on

DStreams?

Some of the output operations are `print()`, `save()` etc.. The `save` operation takes directory to save file into and an optional suffix.

The `print()` takes in the first 10 elements from each batch of the DStream and prints the result.