

Apache Spark

TABLE OF CONTENTS

Preface	2
Introduction to Apache Spark	2
What is Apache Spark?	2
Why Use Apache Spark?	2
Key Features of Apache Spark	2
Spark Components Overview	2
Getting Started with Spark	2
Installation and Setup	2
Using Spark on Local Machine	2
Initializing Spark	2
Resilient Distributed Datasets (RDDs)	3
Creating RDDs	3
Transformations on RDDs	3
Actions on RDDs	3
RDD Persistence	3
Structured APIs: DataFrames and Datasets	3
Creating DataFrames	3
Basic DataFrame Operations	4
Aggregations and Grouping	4
Working with Datasets	4
Spark SQL	4
Registering and Querying Tables	4
Running SQL Queries	4
DataFrame to RDD Conversion	4
Streaming Processing with Spark	4
DStream Creation	4
Transformations on DStreams	4
Output Operations for DStreams	5
Machine Learning with MLlib	5
MLlib Overview	5
Data Preparation	5
Building and Evaluating Models	5
Graph Processing with GraphX	5
Creating Graphs	5
Vertex and Edge RDDs	6
Graph Algorithms	6

Cluster Computing and Deployment	6
Cluster Manager Selection	6
Deploying Spark on Clusters	6
Performance Tuning and Optimization	6
Memory Management	6
Parallelism and Partitions	6
Caching Strategies	6
Interacting with External Data Sources	7
Reading and Writing Data	7
Supported File Formats	7
Connecting to Databases	7
Monitoring and Debugging	7
Spark UI	7
Logging and Debugging	7
Integration with Other Tools	7
Spark and Hadoop	7
Spark and Apache Kafka	7
Spark and Jupyter Notebooks	8
Commonly Used Libraries with Spark	8

PREFACE

This cheatsheet is designed to provide quick access to the most commonly used Spark components, methods, and practices. Whether you're diving into Spark's resilient distributed datasets (RDDs), exploring the DataFrame and SQL capabilities, or harnessing the advanced machine learning libraries through MLlib, this cheatsheet offers bite-sized code snippets and explanations to facilitate your learning.

INTRODUCTION TO APACHE SPARK

WHAT IS APACHE SPARK?

Apache Spark is an open-source, distributed computing system designed for big data processing. It provides an interface for programming entire clusters with implicit data parallelism and fault tolerance. Spark's core abstraction is the Resilient Distributed Dataset (RDD), a fault-tolerant collection of elements that can be processed in parallel.

WHY USE APACHE SPARK?

Spark offers significant advantages over traditional MapReduce-based systems, including faster processing speed due to in-memory computation, a wide range of libraries for various data processing tasks, and support for multiple languages such as Java, Scala, Python, and R.

KEY FEATURES OF APACHE SPARK

- **Speed:** Spark's in-memory processing capability results in faster data processing.
- **Ease of Use:** Provides high-level APIs in languages like Scala, Python, and Java.
- **Versatility:** Supports batch processing, interactive queries, streaming, machine learning, and graph processing.
- **Fault Tolerance:** Recovers lost data using lineage information.
- **Advanced Analytics:** Offers libraries for machine learning (MLlib), graph processing (GraphX), and more.
- **Integration:** Seamlessly integrates with Hadoop, HDFS, and other data sources.

SPARK COMPONENTS OVERVIEW

- **Spark Core:** Foundation of Spark, providing basic functionality like task scheduling, memory management, and fault recovery.
- **Spark SQL:** Enables SQL querying and DataFrame API for structured data processing.
- **Spark Streaming:** Enables processing of real-time data streams.
- **MLlib:** Library for machine learning tasks.
- **GraphX:** Library for graph computation.
- **Cluster Managers:** Supports various cluster managers like Apache Mesos, Hadoop YARN, and Kubernetes.

GETTING STARTED WITH SPARK

INSTALLATION AND SETUP

Apache Spark can be installed on various platforms. Here's a basic guide for setting it up on a local machine

USING SPARK ON LOCAL MACHINE

- Download the latest Spark version from the official website.
- Extract the downloaded archive.
- Set up environment variables, such as `SPARK_HOME` and `PATH`.
- Configure `spark-defaults.conf` for basic settings.

INITIALIZING SPARK

To use Spark in your application, initialize a `SparkSession`

```
import
org.apache.spark.sql.SparkSession;

public class SparkApp {
    public static void main(String[]
args) {
        SparkSession spark =
SparkSession.builder()
            .appName("SparkApp")
```

```

        .master("local[*]") //
Use all available cores
        .getOrCreate();

        // Your Spark application
code here

        spark.stop(); // Stop the
SparkSession
    }
}

```

RESILIENT DISTRIBUTED DATASETS (RDDS)

CREATING RDDS

You can create RDDs from existing data or by parallelizing a collection

```

import
org.apache.spark.api.java.JavaRDD;
import org.apache.spark.SparkConf;
import
org.apache.spark.SparkContext;

SparkConf conf = new
SparkConf().setAppName("RDDExample")
.setMaster("local[*]");
SparkContext sc = new
SparkContext(conf);

List<Integer> data =
Arrays.asList(1, 2, 3, 4, 5);
JavaRDD<Integer> rdd =
sc.parallelize(data);

```

TRANSFORMATIONS ON RDDS

Transformations create a new RDD from an existing one

```

JavaRDD<Integer> squaredRDD =
rdd.map(x -> x * x);
JavaRDD<Integer> filteredRDD =
rdd.filter(x -> x % 2 == 0);

```

```

JavaRDD<Integer> unionRDD =
rdd1.union(rdd2);

```

ACTIONS ON RDDS

Actions return values to the driver program or write data to an external storage system

```

long count = rdd.count();
int firstElement = rdd.first();
List<Integer> collectedData =
rdd.collect();
rdd.saveAsTextFile("output.txt");

```

RDD PERSISTENCE

Caching RDDs in memory can speed up iterative algorithms

```

rdd.persist(StorageLevel.MEMORY_ONLY
());
rdd.unpersist(); // Remove from
memory

```

STRUCTURED APIS: DATAFRAMES AND DATASETS

CREATING DATAFRAMES

DataFrames can be created from various data sources

```

import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import
org.apache.spark.sql.SparkSession;

SparkSession spark =
SparkSession.builder()
.appName("DataFrameExample")
.master("local[*]")
.getOrCreate();

Dataset<Row> df =
spark.read().json("data.json");

```

BASIC DATAFRAME OPERATIONS

Perform various operations on DataFrames

```
df.show();
df.printSchema();
df.select("name").show();
df.filter(df.col("age").gt(21)).show();
df.groupBy("age").count().show();
```

AGGREGATIONS AND GROUPING

Perform aggregations on DataFrames

```
df.groupBy("age").agg(functions.avg("salary"),
functions.max("bonus")).show();
```

WORKING WITH DATASETS

Datasets offer strongly-typed, object-oriented programming interfaces

```
Dataset<Person> people =
df.as(Encoders.bean(Person.class));
people.filter(person ->
person.getAge() > 25).show();
```

SPARK SQL

REGISTERING AND QUERYING TABLES

Register DataFrames as temporary tables for SQL querying

```
df.createOrReplaceTempView("employees");
```

RUNNING SQL QUERIES

Execute SQL queries on registered tables

```
Dataset<Row> results =
```

```
spark.sql("SELECT name, age FROM
employees WHERE age > 25");
results.show();
```

DATAFRAME TO RDD CONVERSION

Convert DataFrames to RDDs when needed

```
JavaRDD<Row> rddFromDF =
df.rdd().toJavaRDD();
```

STREAMING PROCESSING WITH SPARK

DSTREAM CREATION

Create a DStream for streaming processing

```
import
org.apache.spark.streaming.Durations
;
import
org.apache.spark.streaming.api.java.
JavaStreamingContext;

JavaStreamingContext
streamingContext = new
JavaStreamingContext(sparkConf,
Durations.seconds(1));
JavaReceiverInputDStream<String>
lines =
streamingContext.socketTextStream("localhost", 9999);
```

TRANSFORMATIONS ON DSTREAMS

Perform transformations on DStreams

```
JavaDStream<String> words =
lines.flatMap(x ->
Arrays.asList(x.split("
")).iterator());
JavaPairDStream<String, Integer>
wordCounts = words.mapToPair(s ->
new Tuple2<>(s, 1)).reduceByKey((a,
```

```
b) -> a + b);
```

OUTPUT OPERATIONS FOR DSTREAMS

Perform output operations on DStreams

```
wordCounts.print();
wordCounts.saveAsTextFiles("wordcount", "txt");
```

MACHINE LEARNING WITH MLLIB

MLLIB OVERVIEW

MLlib is a powerful library for machine learning tasks

```
import org.apache.spark.ml.Pipeline;
import
org.apache.spark.ml.classification.LogisticRegression;
import
org.apache.spark.ml.evaluation.BinaryClassificationEvaluator;
import
org.apache.spark.ml.feature.VectorAssembler;
import
org.apache.spark.ml.feature.StringIndexer;
```

DATA PREPARATION

Prepare data for machine learning

```
Dataset<Row> rawData =
spark.read().csv("data.csv");
VectorAssembler assembler = new
VectorAssembler()
    .setInputCols(new
String[]{"feature1", "feature2"})
    .setOutputCol("features");
Dataset<Row> assembledData =
assembler.transform(rawData);
```

BUILDING AND EVALUATING MODELS

Build and evaluate a machine learning model

```
StringIndexer labelIndexer = new
StringIndexer()
    .setInputCol("label")
    .setOutputCol("indexedLabel");
LogisticRegression lr = new
LogisticRegression()
    .setMaxIter(10)
    .setRegParam(0.01);
Pipeline pipeline = new Pipeline()
    .setStages(new
PipelineStage[]{labelIndexer,
assembler, lr});
PipelineModel model =
pipeline.fit(trainingData);
Dataset<Row> predictions =
model.transform(testData);
BinaryClassificationEvaluator
evaluator = new
BinaryClassificationEvaluator()
    .setLabelCol("indexedLabel")

    .setRawPredictionCol("rawPrediction");
double accuracy =
evaluator.evaluate(predictions);
```

GRAPH PROCESSING WITH GRAPHX

CREATING GRAPHS

Create a graph in GraphX

```
import
org.apache.spark.graphx.Graph;
import
org.apache.spark.graphx.VertexRDD;
import
org.apache.spark.graphx.util.GraphGenerators;

Graph<Object, Object> graph =
GraphGenerators.logNormalGraph(spark
Context, numVertices, numEPart, mu,
```

```
sigma);
```

VERTEX AND EDGE RDDS

Access vertex and edge RDDs

```
VertexRDD<Object> vertices =
graph.vertices();
EdgeRDD<Object> edges =
graph.edges();
```

GRAPH ALGORITHMS

Apply graph algorithms on the graph

```
import
org.apache.spark.graphx.lib.PageRank
;

Graph<Object, Object> pageRankGraph
=
PageRank.runUntilConvergence(graph,
tolerance);
```

CLUSTER COMPUTING AND DEPLOYMENT

CLUSTER MANAGER SELECTION

Choose a cluster manager for Spark deployment

```
// Set Spark to run on Mesos
SparkConf conf = new SparkConf()
    .setMaster("mesos://mesos-
master:5050")
    .setAppName("SparkApp");

// Set Spark to run on YARN
SparkConf conf = new SparkConf()
    .setMaster("yarn")
    .setAppName("SparkApp");
```

DEPLOYING SPARK ON CLUSTERS

Submit Spark applications to the cluster

```
// Submit using spark-submit script
$ spark-submit --class
com.example.SparkApp --master yarn
--deploy-mode cluster myApp.jar
```

PERFORMANCE TUNING AND OPTIMIZATION

MEMORY MANAGEMENT

Optimize memory usage in Spark

```
// Set memory configurations
conf.set("spark.driver.memory",
"2g");
conf.set("spark.executor.memory",
"4g");

// Enable off-heap memory
conf.set("spark.memory.offHeap.enabl
ed", "true");
conf.set("spark.memory.offHeap.size"
, "2g");
```

PARALLELISM AND PARTITIONS

Adjust parallelism and partitions for better performance

```
// Set the number of executor cores
conf.set("spark.executor.cores",
"4");

// Repartition RDDs for balanced
workloads
JavaRDD<Integer> repartitionedRDD =
rdd.repartition(10);
```

CACHING STRATEGIES

Cache RDDs and DataFrames for repeated computations

```
rdd.persist(StorageLevel.MEMORY_AND_DISK());
df.cache();
```

INTERACTING WITH EXTERNAL DATA SOURCES

READING AND WRITING DATA

Read and write data from/to external sources

```
Dataset<Row> csvData =
spark.read().csv("data.csv");
csvData.write().parquet("data.parquet");
```

SUPPORTED FILE FORMATS

Spark supports various file formats

```
Dataset<Row> parquetData =
spark.read().parquet("data.parquet");
```

CONNECTING TO DATABASES

Connect to databases using JDBC

```
Dataset<Row> jdbcData = spark.read()
    .format("jdbc")
    .option("url",
"jdbc:mysql://host:port/database")
    .option("dbtable", "table")
    .option("user", "username")
    .option("password", "password")
    .load();
```

MONITORING AND DEBUGGING

SPARK UI

Monitor application progress using the Spark UI

```
// Access the Spark UI from the
driver program's URL
http://&#47;&#47;driver-node:4040
```

LOGGING AND DEBUGGING

Use logging for debugging

```
import org.apache.log4j.Logger;
import org.apache.log4j.Level;

Logger.getLogger("org").setLevel(Level.ERROR);
```

INTEGRATION WITH OTHER TOOLS

SPARK AND HADOOP

Spark can work seamlessly with Hadoop

```
// Use HDFS file paths
JavaRDD<String> lines =
sparkContext.textFile("hdfs://namenode:8020/input.txt");
```

SPARK AND APACHE KAFKA

Integrate Spark with Kafka for real-time data processing

```
import
org.apache.spark.streaming.kafka010.
KafkaUtils;
import
org.apache.spark.streaming.kafka010.
LocationStrategies;
import
org.apache.spark.streaming.kafka010.
ConsumerStrategies;

JavaInputDStream<ConsumerRecord<String, String>> kafkaStream =
KafkaUtils.createDirectStream(
    streamingContext,
```



```
LocationStrategies.PreferConsistent(
),

ConsumerStrategies.Subscribe(topics,
kafkaParams)
);
```

SPARK AND JUPYTER NOTEBOOKS

Use [Jupyter Notebooks](#) for interactive data exploration with Spark

```
# Use PySpark in Jupyter Notebook
from pyspark.sql import SparkSession
spark =
SparkSession.builder.appName("SparkA
pp").getOrCreate()
```

COMMONLY USED LIBRARIES WITH SPARK

Library	Description
Spark NLP	Natural Language Processing library for Spark.
Spark Cassandra Connector	Interact with Apache Cassandra.
Spark BigDL	Distributed deep learning library for Spark.
Spark GATK	Genome Analysis Toolkit library for Spark.
Spark TensorFrames	Library for TensorFlow integration with Spark.



JCG delivers over 1 million pages each month to more than 700K software developers, architects and decision makers. JCG offers something for everyone, including news, tutorials, cheat sheets, research guides, feature articles, source code and more.

Copyright © 2014 Exelixis Media P.C. All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by means electronic, mechanical, photocopying, or otherwise, without prior written permission of the publisher.

CHEATSHEET FEEDBACK
WELCOME
support@javacodegeeks.com

SPONSORSHIP
OPPORTUNITIES
sales@javacodegeeks.com