

# **Apache Spark API By Example**

## **A Command Reference for Beginners**

Matthias Langer, Zhen He

Department of Computer Science and Computer Engineering  
La Trobe University  
Bundoora, VIC 3086  
Australia

[m.langer@latrobe.edu.au](mailto:m.langer@latrobe.edu.au), [z.he@latrobe.edu.au](mailto:z.he@latrobe.edu.au)

May 31, 2014

# Contents

<b>1 Preface</b>	<b>5</b>
<b>2 Shell Configuration</b>	<b>6</b>
2.1 Adjusting the amount of memory . . . . .	6
2.2 Adjusting the number of worker threads . . . . .	6
2.3 Adding a Listener to the Logging System . . . . .	6
<b>3 The RDD API</b>	<b>7</b>
3.1 aggregate . . . . .	8
3.2 cartesian . . . . .	10
3.3 checkpoint . . . . .	10
3.4 coalesce, repartition . . . . .	11
3.5 cogroup <sup>[Pair]</sup> , groupWith <sup>[Pair]</sup> . . . . .	11
3.6 collect, toArray . . . . .	12
3.7 collectAsMap <sup>[Pair]</sup> . . . . .	12
3.8 combineByKey <sup>[Pair]</sup> . . . . .	13
3.9 compute . . . . .	13
3.10 context, sparkContext . . . . .	14
3.11 count . . . . .	14
3.12 countApprox . . . . .	14
3.13 countByKey <sup>[Pair]</sup> . . . . .	14
3.14 countByKeyApprox <sup>[Pair]</sup> . . . . .	15
3.15 countByValue . . . . .	15
3.16 countByValueApprox . . . . .	15
3.17 countApproxDistinct . . . . .	16
3.18 countApproxDistinctByKey <sup>[Pair]</sup> . . . . .	16
3.19 dependencies . . . . .	17
3.20 distinct . . . . .	17
3.21 first . . . . .	18
3.22 filter . . . . .	18
3.23 filterWith . . . . .	19
3.24 flatMap . . . . .	20
3.25 flatMapValues <sup>[Pair]</sup> . . . . .	20
3.26 flatMapWith . . . . .	21
3.27 fold . . . . .	21
3.28 foldByKey <sup>[Pair]</sup> . . . . .	22

3.29	foreach . . . . .	22
3.30	foreachPartition . . . . .	22
3.31	foreachWith . . . . .	23
3.32	generator, setGenerator . . . . .	23
3.33	getCheckpointFile . . . . .	23
3.34	preferredLocations . . . . .	24
3.35	getStorageLevel . . . . .	24
3.36	glom . . . . .	25
3.37	groupBy . . . . .	25
3.38	groupByKey <sup>[Pair]</sup> . . . . .	26
3.39	histogram <sup>[Double]</sup> . . . . .	27
3.40	id . . . . .	27
3.41	isCheckpointed . . . . .	28
3.42	iterator . . . . .	28
3.43	join <sup>[Pair]</sup> . . . . .	28
3.44	keyBy . . . . .	29
3.45	keys <sup>[Pair]</sup> . . . . .	29
3.46	leftOuterJoin <sup>[Pair]</sup> . . . . .	30
3.47	lookup <sup>[Pair]</sup> . . . . .	30
3.48	map . . . . .	31
3.49	mapPartitions . . . . .	31
3.50	mapPartitionsWithContext . . . . .	32
3.51	mapPartitionsWithIndex . . . . .	33
3.52	mapPartitionsWithSplit . . . . .	33
3.53	mapValues <sup>[Pair]</sup> . . . . .	34
3.54	mapWith . . . . .	34
3.55	mean <sup>[Double]</sup> , meanApprox <sup>[Double]</sup> . . . . .	35
3.56	name, setName . . . . .	35
3.57	partitionBy <sup>[Pair]</sup> . . . . .	35
3.58	partitioner . . . . .	36
3.59	partitions . . . . .	36
3.60	persist, cache . . . . .	36
3.61	pipe . . . . .	37
3.62	reduce . . . . .	37
3.63	reduceByKey <sup>[Pair]</sup> , reduceByKeyLocally <sup>[Pair]</sup> , reduceByKeyToDriver <sup>[Pair]</sup> . . . . .	37
3.64	rightOuterJoin <sup>[Pair]</sup> . . . . .	38
3.65	sample . . . . .	38
3.66	saveAsHadoopFile <sup>[Pair]</sup> , saveAsHadoopDataset <sup>[Pair]</sup> , saveAsNewAPIHadoopFile <sup>[Pair]</sup> . . . . .	39
3.67	saveAsObjectFile . . . . .	39
3.68	saveAsSequenceFile <sup>[SeqFile]</sup> . . . . .	40
3.69	saveAsTextFile . . . . .	40
3.70	stats <sup>[Double]</sup> . . . . .	41
3.71	sortByKey <sup>[Ordered]</sup> . . . . .	42
3.72	stdev <sup>[Double]</sup> , sampleStdev <sup>[Double]</sup> . . . . .	42

3.73 subtract . . . . .	43
3.74 subtractByKey <sup>[Pair]</sup> . . . . .	43
3.75 sum <sup>[Double]</sup> , sumApprox <sup>[Double]</sup> . . . . .	44
3.76 take . . . . .	44
3.77 takeOrdered . . . . .	45
3.78 takeSample . . . . .	45
3.79 toDebugString . . . . .	46
3.80 toJavaRDD . . . . .	46
3.81 top . . . . .	46
3.82 toString . . . . .	47
3.83 union, ++ . . . . .	47
3.84 unpersist . . . . .	47
3.85 values <sup>[Pair]</sup> . . . . .	48
3.86 variance <sup>[Double]</sup> , sampleVariance <sup>[Double]</sup> . . . . .	48
3.87 zip . . . . .	48
3.88 zipPartitions . . . . .	49
<b>4 Further Topics</b>	<b>51</b>
4.1 Reading from HDFS . . . . .	51

# 1 Preface

Spark is an advanced open-source cluster computing system that is capable of handling extremely large data sets. It was first published by ? and its popularity has increased ever since. Due to its real-time properties and efficient usage of resources, Spark has become a very popular alternative to well established computational software for big data.

Spark is still actively being maintained and further developed by its original creators from UC Berkeley. Hence, this command reference and the associated, including the code-snippets and sample outputs outputs shown, should be considered as a overview of the status-quo of this amazing piece of software technology. Specifically, the API examples in this document are for **Spark version 0.9**. However, we do not expect the API to change much in future releases.

This document does not cover any installation or distribution related topics. For installation instructions, please refer to the Apache Spark website.

[Click here to download full PDF material](#)