

	Issue(s)	Solution(s)	Code Sample
1	High processing cycles.	Principle: "Filter-Early, Filter-Often". Use filter() & map() effectively <b>*before*</b> shuffle ops (e.g., reduceByKey() & groupByKey()). Us them <b>*before*</b> & <b>*after*</b> join() ops.	Use filter() transformations to remove unneeded records. Use map() transformations to project only required fields in RDD.
2	Poor performance in associate operations (e.g., sum() & count())	If you group by a key on a partitioned OR distributed dataset exclusively with the aim of aggregating values for each key, then <b>reduceByKey()</b> is <b>*always*</b> better than using groupByKey().	<code>rdd.map(lambda x: (x[0],1)) \ .groupByKey() \ .mapValues(lambda x: sum(x)) \ .collect()</code> <code>rdd.map(lambda x: (x[0],1)) \ .reduceByKey(lambda x,y:x+y) \ .collect()</code> → <b>Poor code</b> → <b>Good code</b>
3	Poor performance in associate operations (e.g., sum() & count())	If you group by a key on a partitioned OR distributed dataset exclusively with the aim of aggregating values for each key but inputs and outputs to your reduce function are different, then <b>combineByKey()</b> is <b>*always*</b> better than using groupByKey(). Internally, combineByKey combines values of a PairRDD partition by applying an aggregate function.	<code>comb_rdd = student_rdd.map(lambda t: (t[0], (t[1], t[2]))) \ .combineByKey(createCombiner, mergeValue, mergeCombiner) \ .map(lambda t: (t[0], t[1][0]/t[1][1]))</code>
4	Poor performance in associate operations (e.g., sum() & count())	Using <b>foldByKey()</b> is <b>*always*</b> better than using groupByKey() when you wish to aggregate values based on keys and wish to perform an associative operation providing a zero value.	<code>val maxByDept = employeeRDD.foldByKey(("dummy",0.0)) (acc,element)=&gt; if(acc._2 &gt; element._2) acc else element)</code>
5	Memory issues on Spark Workers caused mostly owing to excess Network I/O because of poor design choice; i.e., passing too much data to a function.	Use broadcast method to create a <b>broadcast</b> variable(s) which are shared across Workers using an efficient peer-to-peer sharing protocol, based on BitTorrent. 5 benefits: 1) no shuffle, 2) highly scalable peer-to-peer distribution, 3) replicate data once/Worker instead of once/Task, 4) tasks reuse, 5) provides serialized objects, so they are efficiently read always.	<code>massive_list = [...]</code> \ <code>def big_func(x): # function with massive list \ rdd.map(lambda x: big_func(x)).saveAsTextFile</code> → <b>Poor code</b> <code>stations = sc.broadcast(sdata) \ status = sc.textFile('file:///opt/spark/data/mydata/status') \</code> → <b>Good code</b> <code>.map(lambda x: x.split(',')) \</code> <code>.keyBy(lambda x: x[0])</code>
6	Memory exceptions during collecting data owing to poor design choice.	Best practice is: 1) use take(n) & takeSample() <b>*always*</b> if you need to just inspect instead of collect() & take(). 2) if doing ETL, then persist to a filesystem OR a database instead of doing collect() & take().	
7	Lack of parallelism.	At app level OR using spark-defaults.conf is spark.default.parallelism setting. Best practice here is to have this value equal to OR 2x number of cores on each Spark Worker. You will need to tweak for your own environment, to see what works the best.	Formula to use: <code>spark.default.parallelism = spark.executor.instances * spark.executors.cores * 2</code>
8	Spark Executors idle for long periods of time.	Use dynamic allocation which releases back Executors back to the cluster pool if they are idle for a specified period. Should be implemented typically as a system setting to help maximize system resources.	<code># enable DynamicAllocation disabled by default</code> <code>spark.dynamicAllocation.enabled = True \ spark.dynamicAllocation.minExecutors=n \ spark.dynamicAllocation.maxExecutors=n</code> <code># upperbound is by default set to 60 secs</code>
9	Poor partitioning design choice.	<b>Best practice #1:</b> Avoid small files resulting in too many small partitions. This results in massive overhead of spawning these tasks which is greater than processing required to execute the small tasks.	E.g., filter() ops on partitioned RDD may result in some partitions smaller than others. Solution is to filter() ops with a repartition() OR coalesce() and specify a number less than input RDD. NOTE: repartition() can increase OR reduce number of partitions; coalesce() can only reduce the number of partitions.
10	Poor partitioning design choice.	<b>Best practice #2:</b> Avoid exceptionally large partitions. E.g., loading an RDD from a large compressed file using an unsplittable compression format such as Gzip.	4 solutions: 1) Avoid using unsplittable compression, if possible. 2) Uncompress file locally (e.g., /tmp) before loading into RDD. 3) Repartition <b>*before*</b> a large shuffle operation. 4) Repartition immediately <b>*after*</b> the first transformation against the RDD.
11	Poor partitioning design choice.	<b>Best practice #3:</b> Avoid having fewer partitions than Executors.	Recommended practice is to make this an input parameter to your application and then test out with different input values and optimize for performance. No silver bullet; needs to be found with trial & error.
12	Wrong result in Accumulators.	Accumulators typically used for ops purposes. If used in add-in-place ops to calculate results inside a map() ops, then results can be wrong. <b>Best practice</b> is to always use accumulators only within actions computed by Spark driver, e.g., foreach() action.	Using <i>custom</i> accumulators for accumulating vectors as either lists OR dictionaries. <code>from pyspark import AccumulatorParam</code> <code>class VectorAccumulatorParam(AccumulatorParam): ...</code> <code>vector_acc = sc.accumulator([.....], VectorAccumulatorParam())</code>
13	Actions result in frequent re-evaluations.	Caching an RDD persists data in memory; same routine can then reuse it multiple times when subsequent actions are called without requiring any explicit re-evaluation. Cached partitions are stored in Mem of Executor JVMs on Spark Worker nodes.	<code>words = doc.flatMap(lambda x: x.split()) \ .map(lambda x: (x,1)) \ .reduceByKey(lambda x,y:x+y)</code> <code>words.cache() \ words.count() # triggers computation \ words.take(2) # no computation \ words.count() # no computation</code>
14	Actions result in frequent re-evaluations.	persist() more optimized than cache(). Offers additional storage options including MEMORY_AND_DISK, DISK_ONLY, MEMORY_ONLY_SER, MEMORY_AND_DISK_SER, and MEMORY_ONLY (same as cache() method). Additionally, persist can use replication to persist the same partition on more than one node.	<code>words = doc.flatMap(lambda x: x.split()) \ .map(lambda x: (x,1)) \ .reduceByKey(lambda x,y:x+y)</code> <code>words.persist()</code>
15	Stack overflow issues and long lineage leading to long recovery times.	Using checkpointing eliminates the need for Spark to maintain RDD lineage, especially applicable for streaming and/or iterative processing applications. Checkpointing is expensive so happens only after an action is called.	<code>words = doc.flatMap(lambda x: x.split()) \ .map(lambda x: (x,1)) \ .reduceByKey(lambda x,y:x+y)</code> <code>words.checkpoint() \ words.count() \ words.isCheckpointed() \ words.getCheckpointFile()</code>