

CS535 BIG DATA

PART 1. BATCH COMPUTING MODELS FOR BIG DATA ANALYTICS
1. DISTRIBUTED MODEL FOR SCALABLE
BATCH COMPUTING - MAPREDUCE

Sangmi Lee Pallickara
Computer Science, Colorado State University
<http://www.cs.colostate.edu/~cs535>

9/13/2016 CS535 Big Data - Fall 2016 W4.A.1

FAQs

- Questions about PA1
 - Send an email to cs535@cs.colostate.edu
- FAQ for PA1 page is available at:
 - http://www.cs.colostate.edu/~cs535/FAQ_PA1.html

9/13/2016 CS535 Big Data - Fall 2016 W4.A.2

Objectives

- In-Memory Cluster Computing
 - Programming with Spark

9/13/2016 CS535 Big Data - Fall 2016 W4.A.3

In-Memory Cluster Computing: Apache Spark Closures

9/13/2016 CS535 Big Data - Fall 2016 W4.A.4

Understanding closures

- To execute jobs, Spark breaks up the processing of RDD operations into tasks to be executed by an executor
- Prior to execution, Spark computes the task's **closure**
- The closure is those variables and methods that **must** be **visible for the executor** to perform its computations on the RDD
- This closure is serialized and sent to each executor.

9/13/2016 CS535 Big Data - Fall 2016 W4.A.5

Understanding closures

```
1: int counter = 0;
2: JavaRDD<Integer> rdd = sc.parallelize(data);
3:
4: rdd.foreach(x -> counter += x);
5:
6: println("Counter value: " + counter);
```

- **counter** (in line 4) is referenced within the `foreach` function, it's no longer the **counter** (in line 1) on the driver node
- **counter** (in line 1) will still be zero
- In local mode, in some circumstances the `foreach` function will actually execute within the same JVM as the driver
 - **counter** may be actually updated

9/13/2016 CS535 Big Data - Fall 2016 W4.A.6

Solutions?

- **Accumulator** provides a mechanism for safely updating a variable when execution is split up across worker nodes in a cluster
- Closures (e.g. loops or locally defined methods) should not be used to mutate some global state
 - Spark does not define or guarantee the behavior of mutations to objects referenced from outside the closures

9/13/2016 CS535 Big Data - Fall 2016 W4.A.7

Accumulators [1/4]

- variables that are only "added" to through an associative and commutative operation
 - Efficiently supported in parallel
 - Used to implement counters (as in MapReduce) or sums

```
LongAccumulator accum = sc.sc().longAccumulator();

sc.parallelize(Arrays.asList(1, 2, 3, 4)).foreach(x ->
    accum.add(x));
// ...
// 10/09/29 18:41:08 INFO SparkContext: Tasks finished in
0.317106 s

accum.value();
// returns 10
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.8

Accumulators [2/4]

- Spark natively supports accumulators of numeric types, and programmers can add support for new types

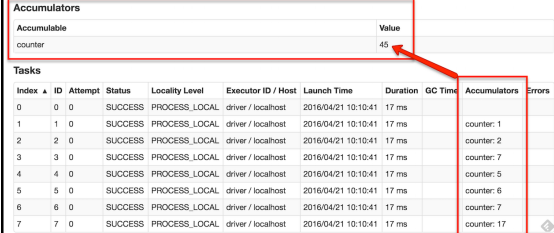
```
class VectorAccumulatorParam implements
AccumulatorParam<Vector> {
    public Vector zero(Vector initialValue) {
        return Vector.zeros(initialValue.size());
    }
    public Vector addInPlace(Vector v1, Vector v2) {
        v1.addInPlace(v2); return v1;
    }
}

// Then, create an Accumulator of this type:
Accumulator<Vector> vecAccum = sc.accumulator(new
Vector(...), new VectorAccumulatorParam());
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.9

Accumulators [3/4]

- If accumulators are created with a name, they will be displayed in Spark's UI



Accumulators		Value
Accumulator	counter	45

9/13/2016 CS535 Big Data - Fall 2016 W4.A.10

Accumulators [4/4]

- Accumulator updates performed inside **actions only**
 - Spark guarantees that each task's update to the accumulator will only be applied once
 - Restarted tasks will not update the value

```
LongAccumulator accum = sc.sc().longAccumulator();
data.map(x -> { accum.add(x); return f(x); });
// Here, accum is still 0 because no actions have caused
the `map` to be computed.
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.11

In-Memory Cluster Computing: Apache Spark

Key-Value pairs

9/13/2016 CS535 Big Data - Fall 2016 W4.A.12

Why Key/Value Pairs?

- Pair RDDs
 - Spark provides special operations on RDDs containing key/value pairs
 - Pair RDDs allow you to act on each key in parallel or regroup data across the network
- `reduceByKey()`
 - Aggregates data separately for each key
- `join()`
 - Merge two RDDs by grouping elements with the same key

9/13/2016 CS535 Big Data - Fall 2016 W4.A.13

Creating Pair RDDs

- Running `map()` function
 - Returns key/value pairs

```
PairFunction < String, String, String > keyData =
  new PairFunction < String, String, String >() {
    public Tuple2 < String, String > call( String x) {
      return new Tuple2(x.split(" ")[ 0], x);
    }
  };

JavaPairRDD <String,String> pairs =
  lines.mapToPair(keyData);
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.14

In-Memory Cluster Computing: Apache Spark

Key-Value pairs

Transformations on Pair RDDs

9/13/2016 CS535 Big Data - Fall 2016 W4.A.15

Transformations on one pair RDD (example: {(1, 2), (3, 4), (3, 6)})

- Pair RDDs are allowed to use all the transformations available to standard RDDs.

Function	purpose	Example	Result
<code>reduceByKey()</code>	Combine values with the same key	<code>rdd.reduceByKey((x, y) => x + y)</code>	{(1,2),(3,10)}
<code>groupByKey()</code>	Group values with the same key	<code>rdd.groupByKey()</code>	{(1,[2]), (3,[4,6])}
<code>combineByKey(createCombiner, mergeValue, mergeCombiners, partitioner)</code>	Combine values with the same key using a different result type		

9/13/2016 CS535 Big Data - Fall 2016 W4.A.16

Transformations on one pair RDD (example: {(1, 2), (3, 4), (3, 6)})

Function	purpose	Example	Result
<code>mapValues(func)</code>	Apply a function to each value of a pair RDD without changing the key	<code>rdd.mapValues(x=> x+1)</code>	{(1,3), (3,5), (3,7)}
<code>flatMapValues(func)</code>	Apply a function that returns an iterator	<code>rdd.flatMapValues(x=>(x to 5))</code>	{(1,2), (1,3), (1,4), (1,5), (3,4), (3,5)}
<code>keys()</code>	Return an RDD of just the keys	<code>rdd.keys()</code>	{1,3,3}
<code>values()</code>	Return an RDD of just the values	<code>rdd.values()</code>	{2,4,6}
<code>sortByKey()</code>	Return an RDD sorted by the key	<code>rdd.sortByKey()</code>	{(1,2), (3,4), (3,5)}

9/13/2016 CS535 Big Data - Fall 2016 W4.A.17

Transformations on one pair RDD (example: {(1, 2), (3, 4), (3, 6)}) other={(3,9)}

Function	purpose	Example	Result
<code>subtractByKey()</code>	Remove elements with a key present in the other RDD	<code>rdd.subtractByKey(other)</code>	{(1,2)}
<code>join()</code>	Inner join	<code>rdd.join(other)</code>	{(3,(4,9)), (3,(6,9))}
<code>rightOuterJoin()</code>	Perform a join where the key must be present in the other RDD	<code>rdd.rightOuterJoin(other)</code>	{(3,(some(4), 9)), (3,(some(6), 9))}
<code>leftOuterJoin()</code>	Perform a join where the key must be present in the first RDD	<code>rdd.leftOuterJoin(other)</code>	{(1,(2, None)), (3,(4, Some(9))), (3,(6, Some(9)))}
<code>coGroup()</code>	Group data from both RDDs sharing the same key	<code>rdd.coGroup(other)</code>	{(1,([2],[])), (3,([4,6],[9]))}

9/13/2016 CS535 Big Data - Fall 2016 W4.A.18

Pair RDDs are still RDDs

- Supports the same functions as RDDs

```
Function < Tuple2 < String, String >, Boolean > longWordFilter =
  new Function < Tuple2 < String, String >, Boolean >() {
    public Boolean call( Tuple2 < String, String > keyValue) {
      return (keyValue._2(). length() < 20);
    }
  };

JavaPairRDD < String, String > result =
  pairs.filter(longWordFilter);
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.19

Filter on Pair RDDs

9/13/2016 CS535 Big Data - Fall 2016 W4.A.20

In-Memory Cluster Computing: Apache Spark

Key-Value pairs Aggregations

9/13/2016 CS535 Big Data - Fall 2016 W4.A.21

Aggregations with Pair RDDs

- Aggregate statistics across all elements with the same key
- reduceByKey()**
 - Similar to `reduce()`
 - Takes a function and use it to combine values
 - Runs several parallel reduce operations
 - One for each key in the dataset
 - Each operations combines values that have the same keys
- `reduceByKey()` is not implemented as an action that returns a value to the user program
 - It returns a new RDD consisting of each key and the reduced value for that key

9/13/2016 CS535 Big Data - Fall 2016 W4.A.22

Example

- Key-value pairs are represented using the `scala.Tuple2` class

```
JavaRDD<String> lines = sc.textFile("data.txt");

JavaPairRDD<String, Integer> pairs =
  lines.mapToPair(s -> new Tuple2(s, 1));

JavaPairRDD<String, Integer> counts =
  pairs.reduceByKey((a, b) -> a + b);
```

- How many times does each line of text occur in a file?

9/13/2016 CS535 Big Data - Fall 2016 W4.A.23

Word count example

```
JavaRDD <String> input = sc.textFile("s3://...")
JavaRDD <String> words = input.flatMap(new FlatMapFunction
< String, String >() {
  public Iterable < String > call( String x) {
    return Arrays.asList( x.split(" "));
  }
});

JavaPairRDD <String,Integer> result = words
  .mapToPair(new PairFunction<String,String,Integer>() {
    public Tuple2 <String,Integer> call(String x) {
      return new Tuple2(x, 1);
    }
  })
  .reduceByKey(new Function2 <Integer,Integer,Integer>(){
    public Integer call(Integer a,Integer b) {
      return a + b;
    }
  });
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.24

combineByKey()

- The most general of the per-key aggregation functions
 - Most of the other per-key combiners are implemented using it
- Allows the user to return values that are not the same type as the input data
- `createCombiner()`
 - If `combineByKey()` finds a new key
 - This happens the first time a key is found in each **partition**, rather than only the first time the key is found in the RDD
- `mergeValue()`
 - If it is not a new value in that partition
- `mergeCombiners()`
 - Merging the results from each partition

9/13/2016 CS535 Big Data - Fall 2016 W4.A.25

Per-key average using `combineByKey()` [1/2]

```
public static class AvgCount implements Serializable {
    public AvgCount( int total, int num) {
        total_ = total;
        num_ = num;
    }
    public int total_;
    public int num_;
    public float avg() {
        return total_ / (float) num_;
    }
}

Function2 < Integer, AvgCount, AvgCount > createAcc = new
Function2 < Integer, AvgCount, AvgCount >() {
    public AvgCount call( Integer x) {
        return new AvgCount( x, 1);
    }
};

Function2 < AvgCount, Integer, AvgCount > addAndCount = new
Function2 < AvgCount, Integer, AvgCount >() {
    public AvgCount call( AvgCount a, Integer x) {
        a.total_ += x;
        a.num_ += 1;
    }
};
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.26

Per-key average using `combineByKey()` [2/2]

```
return a; }
};
Function2 < AvgCount, AvgCount, AvgCount > combine =
new Function2 < AvgCount, AvgCount, AvgCount >() {
    public AvgCount call( AvgCount a, AvgCount b) {
        a.total_ += b.total_;
        a.num_ += b.num_;
        return a;
    }
};

AvgCount initial = new AvgCount( 0,0);
JavaPairRDD < String, AvgCount > avgCounts =
nums.combineByKey( createAcc, addAndCount, combine);
Map < String, AvgCount > countMap = avgCounts.collectAsMap();
for (Entry < String, AvgCount > entry : countMap.entrySet()) {
    System.out.println( entry.getKey() + ":"
        + entry.getValue(). avg());
}
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.27

Tuning the level of parallelism

- When performing aggregations or grouping operations, we can ask Spark to use a specific number of partitions
- `reduceByKey((x, y) => x + y, 10)`
- `repartition()`
 - Shuffles the data across the network to create a new set of partitions
 - Expensive operation
 - Optimized version: `coalesce()`
 - Avoids data movement

9/13/2016 CS535 Big Data - Fall 2016 W4.A.28

groupByKey()

- Group our data using the key in our RDD
- On an RDD consisting of keys of type `K` and values of type `V`
 - Results will be RDD of type `[K, Iterable[V]]`
- `cogroup()`
 - Grouping data from multiple RDDs
 - Over two RDDs sharing the same key type, `K`, with the respective value types `V` and `W` gives us back `RDD[(K, (Iterable[V], Iterable[W]))]`

9/13/2016 CS535 Big Data - Fall 2016 W4.A.29

joins

- Inner join
 - Only keys that are present in both pair RDDs are output
- `leftOuterJoin(other)` and `rightOuterJoin(other)`
 - One of the pair RDDs can be missing the key
- `leftOuterJoin(other)`
 - The resulting pair RDD has entries for each key in the source RDD
- `rightOuterJoin(other)`
 - The resulting pair RDD has entries for each key in the other RDD

9/13/2016 CS535 Big Data - Fall 2016 W4.A.30

In-Memory Cluster Computing: Apache Spark

Key-Value pairs

Actions available on Pair RDDs

9/13/2016 CS535 Big Data - Fall 2016 W4.A.31

Actions on pair RDDs(example({(1,2),(3,4),(3,6)}))

Function	Description	Example	Result
<code>countByKey()</code>	Count the number of elements for each key	<code>rdd.countByKey()</code>	<code>{(1,1), (3,2)}</code>
<code>collectAsMap()</code>	Collect the result as a map to provide easy lookup	<code>rdd.collectAsMap()</code>	<code>Map{(1,2), (3,4), (3,6)}</code>
<code>lookup(key)</code>	Return all values associated with the provided key	<code>rdd.lookup(3)</code>	<code>[4, 6]</code>

9/13/2016 CS535 Big Data - Fall 2016 W4.A.32

In-Memory Cluster Computing: Apache Spark

Data Partitioning

9/13/2016 CS535 Big Data - Fall 2016 W4.A.33

Why partitioning?

- Consider an application that keeps a large table of user information in memory
 - An RDD of (UserID, UserInfo) pairs
 - The application periodically combines this table with a smaller file representing events that happened in the last five minutes

9/13/2016 CS535 Big Data - Fall 2016 W4.A.34

Using `partitionBy()`

- Transforms `userData` to hash-partitioned RDD

9/13/2016 CS535 Big Data - Fall 2016 W4.A.35

PageRank

- Initialize each page's rank to 1.0
- On each iteration, have page p send a contribution of $\text{rank}(p) / \text{numNeighbors}(p)$ to its neighbors (the pages it has link to)
- Set each page's rank to $0.15 + 0.85 * \text{contributionReceived}$

9/13/2016 CS535 Big Data - Fall 2016 W4.A.36

PageRank in Spark (Scala)

```
// Assume that our neighbor list was saved as a Spark
objectFile val links =
  sc.objectFile[(String, Seq[String])](“links”)
    .partitionBy(new HashPartitioner(100))
    .persist()
// Initialize each page's rank to 1.0; since we use
// mapValues, the resulting RDD will have the same
// partitioner as links
val links = // RDD of (url, neighbors) pairs
var ranks = // RDD of (url, rank) pairs
for (i <- 1 to ITERATIONS) {
  val contribs = links.join(ranks)
  .flatMap {
    case (url, (links, rank)) =>
      links.map(dest => (dest, rank/links.size))
  }
  ranks = contribs.reduceByKey(_ + _).mapValues(0.15 +
    0.85 * _)
}
ranks.saveAsTextFile(...)
```

9/13/2016 CS535 Big Data - Fall 2016 W4.A.37

Distributed File Systems GFS

9/8/2015 CS535 Big Data - Fall 2015 W3.A.38

This material is built based on

- Sanjay Ghemawat, Howard Gobioff, Shun-Tak Leung:
The Google file system. Proceedings of SOSP 2003: 29-43
- Andrew Fikes, Storage Architecture and Challenges, Faculty Summit, 2010
 - http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en/us/university/relations/facultysummit2010/storage_architecture_and_challenges.pdf
- Jeff Dean's SOCC keynote, Building Large-Scale Internet Services
 - <http://static.googleusercontent.com/media/research.google.com/en/people/jeff/SOCC2010-keynote-slides.pdf>
- <http://sysmagazine.com/posts/206986/>
- Erasure Coding: Backblaze Open sources Reed-Solomon
 - <https://www.backblaze.com/blog/reed-solomon/>
- An introduction to Reed-Solomon codes
 - http://www.cs.cmu.edu/~guyb/realworld/reedsolomon/reed_solomon_codes.html

9/8/2015 CS535 Big Data - Fall 2015 W3.A.39

Analytics in the Data Center

9/8/2015 CS535 Big Data - Fall 2015 W3.A.40

The Machinery

Servers

- CPUs
- DRAMS
- Disks

Racks

- 40-80 servers
- Ethernet switch

Cluster

- >10,000 nodes

9/8/2015 CS535 Big Data - Fall 2015 W3.A.41

Google Cluster Software Environment

- Clusters contain 1000s of machines, typically one or handful of configurations
 - File system (GFS or Colossus) + cluster scheduling system are core services
- Typically 100s to 1000s of active jobs
 - mix of batch and low-latency, user-facing production jobs

9/8/2015 CS535 Big Data - Fall 2015 W3.A.42				
MapReduce Usage statistics in Google				
	Aug. 04	Mar. 06	Sep.07	May. 10
Number of jobs	29K	171K	2,217K	4,474K
Average completion time (secs)	634	874	395	748
Machine years used	217	2,002	11,081	39,121
Input data read (TB)	3,288	52,254	403,152	946,460
Intermediate data (TB)	758	6,743	34,774	132,960
Output data written (TB)	193	2,970	14,018	45,720
Average worker machines	157	268	394	368

9/8/2015 CS535 Big Data - Fall 2015 W3.A.43				
The Realistic View of a Data Center				
<ul style="list-style-type: none"> Typical first year for a new cluster: <ul style="list-style-type: none"> ~1 network rewiring (rolling downtimes: ~5% of machines over 2-day span) ~20 rack failures (40-80 machines instantly disappear, 1-6 hours to get back) ~5 racks go wonky (40-80 machines see 50% packet loss) ~8 network maintenances (4 might cause ~30-minute random connectivity losses) ~12 router reloads (takes out DNS and external IPs for a couple minutes) ~3 router failures (have to immediately pull traffic for an hour) ~dozens of minor 30-second blips for DNS ~1000 individual machine failures ~thousands of hard drive failures <ul style="list-style-type: none"> slow disks, bad memory, misconfigured machines, flaky machines, etc. Long distance links <ul style="list-style-type: none"> Reliability/availability must come from software 				

9/8/2015 CS535 Big Data - Fall 2015 W3.A.44				
Numbers we should know [1/2]				
<ul style="list-style-type: none"> Level 1 cache reference <ul style="list-style-type: none"> 0.5 ns Branch misprediction <ul style="list-style-type: none"> 5 ns Level 2 cache reference <ul style="list-style-type: none"> 7 ns Mutex lock/unlock <ul style="list-style-type: none"> 25 ns Main memory reference <ul style="list-style-type: none"> 100 ns Compress 1KB with cheap compression algorithm <ul style="list-style-type: none"> 3,000 ns 				

9/8/2015 CS535 Big Data - Fall 2015 W3.A.45				
Numbers we should know [2/2]				
<ul style="list-style-type: none"> Read 1 MB sequentially from memory <ul style="list-style-type: none"> 250,000 ns Round trip within the same datacenter <ul style="list-style-type: none"> 500,000 ns Disk seek <ul style="list-style-type: none"> 10,000,000 ns Read 1 MB sequentially from disk <ul style="list-style-type: none"> 20,000,000 ns Send packet CA->Netherlands->CA <ul style="list-style-type: none"> 150,000,000 ns 				

9/8/2015 CS535 Big Data - Fall 2015 W3.A.46				
Back of the Envelope Calculation				
<ul style="list-style-type: none"> How long to generate an image results page (30 thumbnails)? Design 1: Read serially, thumbnail images (256KB) on the fly <ul style="list-style-type: none"> 30 seeks * 10 ms/seek + 30 * 256K / 30 MB/s = 560 ms Design 2: Issue reads in parallel: <ul style="list-style-type: none"> 10 ms/seek + 256K read / 30 MB/s = 18 ms Lots of variations: <ul style="list-style-type: none"> caching (single images? whole sets of thumbnails?) pre-computing thumbnails ... 				

9/8/2015 CS535 Big Data - Fall 2015 W3.A.47				
Storage Software: GFS				
<ul style="list-style-type: none"> Google's first cluster-level file system (2003) <ul style="list-style-type: none"> Designed for batch applications with large files Single master for metadata and chunk management Chunks are typically replicated 3x for reliability Lessons <ul style="list-style-type: none"> Scaled to approximately 50M files, and 10PB Large files increased application complexity Not appropriate for latency sensitive applications Scaling limits added management overhead 				

9/8/2015 CS535 Big Data - Fall 2015 W3.A.48

Storage Software: Colossus (GFS2)

- Next-generation cluster-level file system
- Automatically sharded metadata layer
 - Data typically written using Reed-Solomon (1.5x)
 - Client-driven replication, encoding and replication
 - Metadata space has enabled availability
- Why Reed-Solomon?
 - Cost. Especially with cross cluster replication
 - More flexible cost vs. availability choice

9/8/2015 CS535 Big Data - Fall 2015 W3.A.49

Storage Landscape

- Early Google:
 - US-centric traffic
 - Batch, latency-insensitive indexing processes
 - Document "snippets" serving (single seek)
- Current day:
 - World-wide traffic
 - Continuous crawl and indexing processes (Caffeine)
 - Seek-heavy, latency-sensitive apps (Gmail)
 - Person-to-person, person-to-group sharing (Docs)

9/8/2015 CS535 Big Data - Fall 2015 W3.A.50

Storage Landscape: Flash (SSDs)

- Important future directions:
 - More workloads that are increasingly seek heavy
 - 50-150x less expensive than disk per random read
 - Best usage is still being explored
- Concerns:
 - Availability of devices
 - 17-32x more expensive per GB than disk
 - Endurance not yet proven in the field