Indian Institute of Science Bangalore, India भारतीय विज्ञान संस्थान बंगलौर, भारत

DS221 | 19 Sep – 19 Oct, 2017 Data Structures, Algorithms & Data Science Platforms

Yogesh Simmhan

simmhan@cds.iisc.ac.in



©Department of Computational and Data Science, IISc, 2016 This work is licensed under a <u>Creative Commons Attribution 4.0 International License</u> Copyright for external content used with attribution is retained by their original authors





L5: Big Data Platforms

Spark, Storm, Giraph

Slide Credits:

- https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf
- https://www.slideshare.net/deanchen11/scala-bay-spark-talk
- https://databricks-training.s3.amazonaws.com/slides/advanced-spark-training.pdf
- Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing, M. Zaharia, et al., NSDI 2012
- http://spark.apache.org/docs/latest/programming-guide.html

17-0ct-17



What is Big Data?





The term **is** fuzzy ... Handle with care!



Wordle of "Thought Leaders'" definition of Big Data, © Jennifer Dutcher, 2014 https://datascience.berkeley.edu/what-is-big-data/

05-Jan-17



So...What is Big Data?

Data whose characteristics exceeds the capabilities of conventional *algorithms, systems and techniques* to derive useful value.

https://www.oreilly.com/ideas/what-is-big-data





And, where does Big Data come from?



Web & Social Media

Web search, Social Networks & Micro-blogs





Web & Social Media

Social Networks & Micro-blogs



Facebook's mobile ad revenue as a share of total ad revenue



THE WALL STREET JOURNAL.

1.79 billion monthly active users as of September 30, 2016

05-Jan-17

https://www.statista.com/statistics/264810/number-of-monthly-active-facebook-users-worldwide/ http://www.wsj.com/articles/facebook-profit-jumps-sharply-1478117646 http://newsroom.fb.com/company-info/



Amazon Annual Revenue (\$ Billions)

1998-2013 (Actual), 2014-2015 (Projected)

6.9 8.5 10.7

Enterprises & Government

Online retail & eCommerce



Source: CRISIL Research

http://blogs.ft.com/beyond-brics/2014/02/28/online-retail-in-india-learning-to-evolve/

http://www.peridotcapital.com/2014/04/amazon-sales-growth-projections-for-next-two-years-appear-overly-optimistic.html

05-Jan-17

107.6

89.9

74.5



Enterprises & Government: Finance

Mobile Transactions & FinTech

ASIA/PACIFIC (USERS IN MILLIONS)



Since November 8, 2016, Paytm has surpassed its metrics -tripling transactions per day to 7.5 million

http://www.pymnts.com/in-depth/2015/mobile-transactions/ Is Paytm the Xerox of mobile payments?, ETtech.com-03-Jan-2017

Internet of Everything

- Personal Devices
 - Smart Phones, Fitbit
- Smart Appliances
- Smart Cities
 - Power, Water, Transportation, Environment
- Smart Retail
- Millions of sensor data streams



smartx.cds.iisc.ac.in



Why is Big Data Difficult?





05-Jan-17





captures

1 TR OF TRADE



100 SENSORS

that monitor items such as

data encompasses infor

internal and external source







U.S. have at least

100 TERABYTES

[100,000 GIGABYTES] of data stored

Modern cars have close to 100 SENSORS

that monitor items such as fuel level and tire pressure



1 TB OF TRADE

The New York Stock Exchange

captures

By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

 almost 2.5 connections per person on earth

Velocity ANALYSIS OF STREAMING DATA

YYYYYYYYY



and services that the v But what exactly is big massive amounts of dat

As a leader in the sec break big data into for Velocity, Variety and Ver

Depending on the indust data encompasses inf internal and external sou social media, enterpris mobile devices. Compa adapt their products an customer needs, opt infrastructure, and find

By 2015 **4.4 MILLION IT JO**

will be created globally with 1.9 million in the





R V's g

d music downloads to web ecords, data is recorded, to enable the technology world relies on every day, data, and how can these ta be used?

ctor, IBM data scientists our dimensions: **Volume,** r<mark>acity</mark>

stry and organization, big formation from multiple urces such as transactions, As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES

[161 BILLION GIGABYTES]



30 BILLION PIECES OF CONTENT

are shared on Facebook every month

1 IN 3 BUSINESS

don't trust the information

LEADERS



Variety

DIFFERENT

FORMS OF DATA

By 2014, it's anticipated there will be 420 MILLION WEARABLE, WIRELESS HEALTH MONITORS

4 BILLION+ HOURS OF VIDEO

are watched on YouTube each month



400 MILLION TWEETS

are sent per day by about 200 million monthly active users

Poor data quality costs the US economy around



rld relies on every day. ata, and how can these be used?

r, IBM data scientists r dimensions: **Volume, ity**

y and organization, big mation from multiple ces such as transactions, content, sensors and es can leverage data to services to better meet mize operations and ew sources of revenue.

S

o support big data, Inited States





1 IN 3 BUSINESS LEADERS don't trust the information they use to make decisions



Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate

Veracity UNCERTAINTY OF DATA





Data Analysis Lifecycle





Data Platforms

- Acquire, manage, process Big Data
- At large scales
- To meet application needs



Distributed Systems

Distributed Computing

- Clusters of machines
- Connected over network
- Distributed Storage
 - Disks attached to clusters of machines
 - Network Attached Storage
- How can we make effective use of multiple machines?

Commodity clusters vs. HPC clusters

- Commodity: Available off the shelf at large volumes
- Lower Cost of Acquisition
- Cost vs. Performance
 - Low disk bandwidth, and high network latency
 - CPU typically comparable (Xeon vs. i3/5/7)
 - Virtualization overhead on Cloud

How can we use many machines of modest capability?

Growth of Cloud Data Centers

Figure 17. Global Data Center Workloads by Applications





Ideal Strong/Weak Scaling



22 Scaling Theory and Machine Abstractions, Martha A. Kim, October 10, 2012



Scalability

- Strong vs. Weak Scaling
- Strong Scaling: How the performance varies with the # of processors for a *fixed total problem size*
- Weak Scaling: How the performance varies with the # of processors for a *fixed problem size per* processor
 - Big Data platforms are intended for "Weak Scaling"



Ease of Programming

- Programming distributed systems is difficult
 - Divide a job into multiple tasks
 - Understand dependencies between tasks: Control, Data
 - Coordinate and synchronize execution of tasks
 - Pass information between tasks
 - Avoid race conditions, deadlocks
- Parallel and distributed programming models/languages/abstractions/platforms try to make these easy
 - E.g. Assembly programming vs. C++ programming
 - E.g. C++ programming vs. Matlab programming



Availability, Failure

- Commodity clusters have lower reliability
 - Mass-produced
 - Cheaper materials
 - Smaller lifetime (~3 years)
- How can applications easily deal with failures?
- How can we ensure availability in the presence of faults?



Early Technologies

- MapReduce is a distributed data-parallel programming model from Google
- MapReduce works best with a distributed file system, called Google File System (GFS)
- Hadoop is the open source framework implementation from Apache that can execute the MapReduce programming model
- Hadoop Distributed File System (HDFS) is the open source implementation of the GFS design
- Elastic MapReduce (EMR) is Amazon's PaaS



Platforms...Think in terms of Stacks Cloudera

Cloudera's Distribution for Hadoop

UI Framework Hue			S	DK	Hu	e SDK	
Workflow 002	zie	Scheo	duling <i>Oozie</i>		Metada	ata	Hive
Data		Languages	, Compilers _p	Pig/ Hive	Fast		
Integration <i>Flume, Sqoop</i>		Gh a	doop		read/w access	rite	HBase
		Coordi	ination		2	Zookee	eper



Platforms...Think in terms of Stacks BDAS





Platforms...Think in terms of Stacks HortonWorks





Apache Spark

Slides & Additional Reading Courtesy

https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf

Resilient Distributed Datasets, Matei Zaharia

http://spark.apache.org/docs/2.1.1/programming-guide.html

http://spark.apache.org/docs/latest/api/java/index.html https://www.gitbook.com/book/jaceklaskowski/mastering-apache-spark/details Apache Spark Internals, Pietro Michiardi, Eurecom



Why Spark?



- Ease of language definition
 - Typing, dataflows,
 - But Pig, Hive, HBase, etc. give you that
- Better performance using "In memory" compute
 - Multiple stages part of same job
 - Lazy evaluation, caching/persistence



In-memory computation

Operate on data in (distributed) memory

- Allows many operations to be performed locally
- Write to disk only when data sharing required across workers
- This is unlike others like Hadoop Map/Reduce



RDD: The Secret Sauce

- RDD: Resilient Distributed Dataset
 - Immutable, partitioned collection of tuples
 - Operated on by *deterministic* transformations
 - Object-oriented flavor
 - RDD.operation() \rightarrow RDD
- Recovery by re-computation
 - Maintains lineage of transformations
 - Recompute missing partitions if failure happens
 - Not possible/not automatic in Pig
- Allows caching & persistence for reuse







RDD Operations

		flatMan	into Dataflows
Transformations (define a new RDD)	map filter sample groupByKey reduceByKey	union join cogroup cross	
	sortByKey	mapValue	s
Actions (return a result to driver program)	colle redu cou sav looku	ect uce int /e pKey	

Allows

composability

A Sample Spark Program

- Counts the number of bytes in a line, and sums the count per line
- Uses lambda expressions

// Cache RDD in-memory for future use in this app lineLengths.persist(StorageLevel.MEMORY_ONLY());

A Sample Spark Program

Can pass complex functions as well

```
class GetLength implements Function<String, Integer> {
   public Integer call(String s) { return s.length(); }
}
class Sum implements Function2<Integer, Integer, Integer> {
   public Integer call(Integer a, Integer b) { return a + b; }
}
```

```
JavaRDD<String> lines = sc.textFile("data.txt");
JavaRDD<Integer> lineLengths = lines.map(new GetLength());
int totalLength = lineLengths.reduce(new Sum());
```



RDD Partitions

- RDD is internally a collection of partitions
 - Each partition holds a list of items
- Partitions may be present on a different machine
 - Partition is the unit of execution
 - Partition is the unit of parallelism
- They are immutable
 - Each transformation on an RDD generates a new RDD with different partitions
 - Allows recovery of individual partitions

Partition 1 RDD	Partition 2	Partition 3
	Input Format	
Chunk 1 Data in HDFS	Chunk 2	Chunk 3



Creating RDD

- Load external data from distributed storage
- Create logical RDD on which you can operate
- Support for different input formats
 HDFS files, Cassandra, Java serialized, directory, gzipped
- Can control the number of partitions in loaded RDD
 - Default depends on external DFS, e.g. 128MB on HDFS

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



RDD Operations

- Transformations
 - From one RDD to one or more RDDs
 - Lazy evaluation...use with care
 - Executed in a distributed manner
- Actions
 - Perform aggregations on RDD items
 - Return single (or distributed) results to "driver" code
- RDD.collect() brings RDD partitions to single driver machine



Caution: Local Variables & Closures

- Caution: Cannot pass "local" driver variables to lambda expressions/anonymous classes....only final
 - Will fail when distributed

```
int counter = 0;
JavaRDD<Integer> rdd = sc.parallelize(data);
// Wrong: Don't do this!!
rdd.foreach(x -> counter += x);
println("Counter value: " + counter);
```

RDD and PairRDD

- RDD is logically a collection of items with a generic type
- PairRDD is like a "Map", where each item in collection is a <key,value> pair, each a generic type
- Transformation functions use RDD or PairRDD as input/output
- E.g. Map-Reduce

```
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);
```



Transformation	Meaning
map(func)	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
filter(func)	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
flatMap(func)	Similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item).

- JavaRDD<R> map(Function<T,R> f) : 1:1 mapping from input to output. Can be different types.
- JavaRDD<T> filter(Function<T,Boolean> f) : 1:0/1 from input to output, same type.
- JavaRDD<U> flatMap(FlatMapFunction<T,U> f): 1:N mapping from input to output, different types.



mapPartitions(func)

Similar to map, but runs separately on each partition (block) of the RDD, so *func* must be of type Iterator<T> => Iterator<U> when running on an RDD of type T.

- Earlier Map and Filter operate on one item at a time. No state across calls!
- JavaRDD<U> mapPartitions(FlatMapFunc<Iterator<T>,U> f)
- mapPartitions has access to iterator of values in entire partition, jot just a single item at a time.



<pre>sample(withReplacement, fraction, seed)</pre>	Sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator seed.
union(otherDataset)	Return a new dataset that contains the union of the elements in the source dataset and the argument.

- JavaRDD<T> sample(boolean withReplacement, double fraction): fraction between [0,1] without replacement, >0 with replacement
- JavaRDD<T> union(JavaRDD<T> other): Items in other RDD added to this RDD. Same type. Can have duplicate items (i.e. not a 'set' union).

intersection(otherDataset)	Return a new RDD that contains the intersection of elements in the source dataset and the argument.
distinct([numTasks]))	Return a new dataset that contains the distinct elements of the source dataset.

- JavaRDD<T> intersection(JavaRDD<T> other): Does a set intersection of the RDDs. Output will not have duplicates, even if inputs did.
- JavaRDD<T> distinct(): Returns a new RDD with unique elements, eliminating duplicates.



Transformations: **PairRDD**

groupByKey([numTasks])	When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable <v>) pairs. Note: If you are grouping in order to perform an aggregation (such as a sum or average) over each key, using reduceByKey or aggregateByKey will yield much better performance. Note: By default, the level of parallelism in the output depends on the number of partitions of the parent RDD. You can pass an optional numTasks argument to set a different number of tasks.</v>
<pre>reduceByKey(func, [numTasks])</pre>	When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function <i>func</i> , which must be of type $(V,V) => V$. Like in groupBykey, the number of reduce tasks is configurable through an optional second argument.

- JavaPairRDD<K,Iterable<V>> groupByKey(): Groups values for each key into a single iterable.
- JavaPairRDD<K,V> reduceByKey(Function2<V,V,V> func) : Merge the values for each key into a single value using an <u>associative</u> and <u>commutative</u> reduce function. Output value is of same type as input.
- For aggregate that returns a different type?
- numPartitions can be used to generate output RDD with different number of partitions than input RDD.

aggregateByKey(zeroValue)(seqOp, combOp, [numTasks])	When called on a dataset of (K, V) pairs, returns a dataset of (K, U) pairs where the values for each key are aggregated using the given combine functions and a neutral "zero" value. Allows an aggregated value type that is different than the input value type, while avoiding unnecessary allocations. Like in groupBykey, the number of reduce tasks is configurable through an optional second argument.
<pre>sortByKey([ascending], [numTasks])</pre>	When called on a dataset of (K, V) pairs where K implements Ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean ascending argument.

- JavaPairRDD<K,U> aggregateByKey(U zeroValue, Function2<U,V,U> seqFunc, Function2<U,U,U> combFunc) : Aggregate the values of each key, using given combine functions and a neutral "zero value".
 - **SeqOp** for merging a V into a U within a partition
 - **CombOp** for merging two U's, within/across partitions
- JavaPairRDD<K,V> sortByKey(Comparator<K> comp): Global sort of the RDD by key
 - <u>Each partition</u> contains a sorted range, i.e., output RDD is rangepartitioned.
 - Calling collect will return an ordered list of records



join(otherDataset, [numTasks])	When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with
	all pairs of elements for each key. Outer joins are supported through leftouterjoin, rightouterjoin, and fullouterjoin.
cartesian(otherDataset)	When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).

- JavaPairRDD<K, Tuple2<V,W>>
 join(JavaPairRDD<K,W> other, int numParts):
 Matches keys in *this* and *other*. Each output pair is
 (k, (v1, v2)). Performs a hash join across the cluster.
- JavaPairRDD<T,U> cartesian(JavaRDDLike<U,?> other): Cross product of values in each RDD as a pair



Actions

reduce(func)	Aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
collect()	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
count()	Return the number of elements in the dataset.
countByKey()	Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.



Actions

first()	Return the first element of the dataset (similar to take(1)).
take(n)	Return an array with the first <i>n</i> elements of the dataset.
takeSample (<i>withReplacement</i> , <i>num</i> , [seed])	Return an array with a random sample of <i>num</i> elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.



RDD Persistence & Caching

- RDDs can be reused in a dataflow
 - Branch, iteration
- But it will be re-evaluated each time it is reused!
- Explicitly persist RDD to reuse output of a dataflow path multiple times
- Multiple storage levels for persistence
 - Disk or memory
 - Serialized or object form in memory
 - Partial spill-to-disk possible
 - Cache indicates "persist" to memory



RePartitioning

repartition

public JavaRDD<T> repartition(int numPartitions)

Return a new RDD that has exactly numPartitions partitions.

Can increase or decrease the level of parallelism in this RDD. Internally, this uses a shuffle to redistribute data.

If you are decreasing the number of partitions in this RDD, consider using coalesce, which can avoid performing a shuffle.

coalesce

Return a new RDD that is reduced into numPartitions partitions.



From DAG to RDD lineage





Samples: Word Count

rdd = sc.textFile("hdfs://..."); words = rdd.flatMap(x -> x.split(" ")); result = words.map(x->(x,1)). reduceByKey((x, y): x + y);



Samples: Per-key average



key	value
panda	(1, 2)
pink	(7, 2)
pirate	(3, 1)

sumCount =

rdd.mapValues(x -> (x,1)).
reduceByKey((x, y) ->
 (x[0]+y[0], x[1]+y[1]))



Additional Topics



Sample: PageRank

```
// URL neighbor URL
JavaRDD<String> lines =
spark.read().textFile(args[0]).javaRDD();
// Loads all URLs from input file and initialize their
neighbors.
JavaPairRDD<String, Iterable<String>> links =
lines.mapToPair(s -> {
    String[] parts = SPACES.split(s);
    return new Tuple2<>(parts[0], parts[1]);
    }).distinct().groupByKey().cache();
```

// Loads all URLs with other URL(s) link to from input file
and initialize ranks of them to one.
JavaPairRDD<String, Double> ranks = links.mapValues(rs->1.0);



// Calculates and updates URL ranks continuously using PageRank algorithm.
for (int current = 0; current < Integer.parseInt(args[1]); current++) {
 // Calculates URL contributions to the rank of other URLs.
 JavaPairRDD<String, Double> contribs = links.join(ranks).values()
 .flatMapToPair(s -> { // _1 = adj list, _2 = ranks
 int urlCount = Iterables.size(s._1());
 List<Tuple2<String, Double>> results = new ArrayList<>();
 for (String n : s._1) { // Send rank value to neighbor
 results.add(new Tuple2<>(n, s._2() / urlCount));
 }
 return results.iterator();

```
});
```

```
// Re-calculates URL ranks based on neighbor contributions.
ranks = contribs.reduceByKey(new Sum()).mapValues(sum -> 0.15 + sum * 0.85);
// Collects all URL ranks and dump them to console.
List<Tuple2<String, Double>> output = ranks.collect();
for (Tuple2<?,?> tuple : output) {
    System.out.println(tuple._1() + " has rank: " + tuple._2() + ".");
}
```

Distributed Execution

RDD Objects

Scheduler (DAGScheduler)





Execution Dependency

NARROW DEPENDENCY: Each partition of the parent RDD is used by at most one partition of the child RDD. Task can be executed locally and we don't have to shuffle. **WIDE DEPENDENCY**: Multiple child partitions may depend on one partition of the parent RDD. We have to shuffle data *unless the parents are hash-partitioned*





join with co-partitioned inputs



Lazy Execution





Execution Planning: Word Count

- lines = sc.textFile("input")
- words = lines.flatMap(1 -> l.split(" "))
- ones = words.map(w -> 1)
- counts = ones.reduceByKey((a,b) -> a+b)
- result = counts.collectAsMap()

- RDD lineage DAG is built on driver side
- An action triggers the Job Submission

63





Hash Shuffle



- Distributed data from one set of partitions to another, on different machines
- "Keys" used for routing from Map to Reduce
 - Map writes one "bucket" per key. Reduce reads "buckets" with its key.

17-0ct-17



Schedule

- Thu Oct 12, 1130-1 PM guest lecture on Storm and Giraph
 - Programming streaming data & linked data
- Spark Lecture
 - 530-7PM Thu, Oct 12
 - CDS 202
- Mid-term 2: Thu, Oct 19 1130am-1pm (10% weightage)
- Spark Tutorial
 - 530-7PM Tue, Oct 24
 - CDS 309
- Assignment 3 on PySpark to be posted by Oct 20, due by Oct 30 (5% weightage)