Big Data for Scientists

John Urbanic Parallel Computing Scientist Pittsburgh Supercomputing Center

Copyright 2022

Big data is a broad term for data sets so large or complex that traditional data processing applications are inadequate.

-Wikipedia

Once there was only small data...



A classic amount of "small" data

Find a tasty appetizer – Easy!

Find something to use up these oranges – grumble...

What if....



Less sophisticated is sometimes better...



"Chronologically" or "geologically" organized. Familiar to some of you at tax time. Get all articles from 2007.

Get all papers on "fault tolerance" – grumble and cough

Indexing will determine your individual performance. Teamwork can scale that up.



The culmination of centuries...



Find books on Modern Physics (DD# 539)

Find books by Wheeler

where he isn't the first author – grumble...





Then data started to grow.

1956 IBM Model 350



5 MB of data!

But still pricey. \$

Better think about what you want to save.

And finally got **BIG**.

8TB for \$130





*Actually, a silly estimate. The original referen 2013 the digital collection alone was 3PB.

Genome sequencers (Wikipedia Commons)



Science...

Storage got cheap

So why not keep it all?

IoT

Today data is a hot commodity \$

Facebook

And we got better at generating it

Facebook Deep Learning





۱ŀ

yelp&

You Tube

wikipedia 208TB, and in get_involved/blog/bioblitzinsects-reviewed

1/biodiv whales walrus.html

A better sense of biggish

Size

- 1000 Genomes Project
 - AWS hosted
 - 260TB
- Common Crawl
 - Hosted on Bridges
 - 300-800TB+

Throughput

- Square Kilometer Array
 - Building now
 - Exabyte of raw data/day compressed to 10PB
- Internet of Things (IoT) / motes
 - Endless streaming

Records

- GDELT (Global Database of Events, Language, and Tone) (also soon to be hosted on Bridges)
 - Only about 2.5TB per year, but...
 - 250M rows and 59 fields (BigTable)
 - *"during periods with relatively little content, maximal translation accuracy can be achieved, with accuracy linearly degraded as needed to cope with increases in volume in order to ensure that translation always finishes within the 15 minute window.... and prioritizes the highest quality material, accepting that lower-quality material may have a lower-quality translation to stay within the available time window."*

3 V's of Big Data

- Volume
- Velocity
- Variety

Good Ol' SQL couldn't keep up.

SELECT NAME, NUMBER, FROM PHONEBOOK

Why it *wasn't* fashionable:

- Schemas set in stone:
 - Need to define before we can add data
 - Not a fit for agile development
 "What do you mean we didn't plan to keep logs of everyone's heartbeat?"
- Queries often require accessing multiple indexes and joining and sorting multiple tables
- Sharding isn't trivial
- Caching is tough
 - ACID (Atomicity, Consistency, Isolation, Durability) in a *transaction* is costly.





So we gave up: Key-Value

Redis, Memcached, Amazon DynamoDB, Riak, Ehcache

GET foo

- Certainly agile (no schema)
- Certainly scalable (linear in most ways: hardware, storage, cost)
- Good hash might deliver fast lookup
- Sharding, backup, etc. could be simple
- Often used for "session" information: online games, shopping carts

foo	bar
2	fast
6	0
9	0
0	9
text	pic
1055	stuff
bar	foo

GET cart:joe:15~4~7~0723

How does a pile of unorganized data solve our problems?

Sure, giving up ACID buys us a lot performance, but doesn't our crude organization cost us something? Yes, but remember these guys?



This is what they look like today.



Document



GET foo

- Value must be an object the DB can understand
- Common are: XML, JSON, Binary JSON and nested thereof
- This allows server side operations on the data

GET plant=daisy

- Can be quite complex: Linq query, JavaScript function
- Different DB's have different update/staleness paradigms

foo	
2	<pre><catalogs cflant=""></catalogs></pre>
6	JSON
9	XML
0	Binary JSON
bar	JSON XML
12	XML XML

Wide Column Stores

Cassandra Google BigTable

SELECT Name, Occupation FROM People WHERE key IN (199, 200, 207);

- No predefined schema
- Can think of this as a 2-D key-value store: the value may be a key-value store itself
- Different databases aggregate data differently on disk with different optimizations

Кеу			
Joe	Email: joe@gmail	Web: www.joe.com	
Fred	Phone: 412-555-3412	Email: fred@yahoo.com	Address: 200 S. Main Street
Julia	Email: julia@apple.com		
Мас	Phone: 214-555-5847		



- Great for semantic web
- Great for graphs 😕 •
- Can be hard to visualize
- Serialization can be difficult \bullet
- Queries more complicated \bullet



From <u>PDX Graph Meetup</u>

Queries

SPARQL, Cypher

SPARQL (W3C Standard)

- Uses Resource Description Framework format
 - triple store
- RDF Limitations
 - No named graphs
 - No quantifiers or general statements
 - "Every page was created by some author"
 - "Cats meow"
- Requires a schema or *ontology* (RDFS) to define rules
 - "The object of 'homepage' must be a Document."
 - "Link from an actor to a movie must connect an object of type Person to an object of type Movie."



Cypher (Neo4J only)

- No longer proprietary
- Stores whole graph, not just triples
- Allows for named graphs
- …and general Property Graphs (edges and nodes may have values)

```
SMATCH (Jack:Person
```

```
{ name: 'Jack Nicolson' }) - [:ACTED_IN] - (movie:Movie)
RETURN movie
```

Graph Databases

- These are not curiosities, but are behind some of the most high-profile pieces of Web infrastructure.
- They are definitely *big* data.

Microsoft Bing Knowledge Graph	Search and conversations.	~2 billion primary entries ~55 billion facts
Facebook		~50 million primary entries ~500 million assertions
Google Knowledge Graph	Search and conversations.	~1 billion entries ~55 billion facts
LinkedIn graph		590 million members30 million companies

Hadoop & Spark

What kind of databases are they?

Frameworks for Data

These are both frameworks for distributing and retrieving data. Hadoop is focused on disk based data and a basic map-reduce scheme, and Spark evolves that in several directions that we will get in to. Both can accommodate multiple types of databases and *achieve their performance gains by using parallel workers*.



The mother of Hadoop was necessity. It is trendy to ridicule its primitive design, but it was the first step.



What exactly is this Hadoop "framework"?

- Programming platform
- Distributed filesystem
- Parallel execution environment
- Software ecosystem

Programming = MapReduce



Transponder ID -> Geo Coordinates 00154301 -> 59.33, 177.60 04435354 -> 56.71, 171.73 04539340 -> 25.18, -118.89 Only keep data for Bering Sea 00154301 -> 59.33, 177.60 04435354 -> 56.71, 171.73 Find biggest change at each transponder in last 24h Keep any over 20 degrees

00154301 -> 30

00154301 -> 30 04435354 -> 5

HDFS: Hadoop Distributed File System

- Replication
 - Failsafe
 - Predistribution
- Write Once Read Many (WORM)
 - No Random Access (contrast with RDBMS)
- Requires underlying filesystem

YARN



Source: https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/YARN.html

Using Hadoop

- 1. Load data to HDFS
- 2. Write a program
 - Java (compile/jar/run)
 - Hadoop Streaming
 - Mapper and reducer scripts read/write stdin/stdout
 - Use built-in utilities (wc, grep, cat)
 - Write in any language (python)
- 3. Submit a job

Shakespeare Using Hadoop

(from the future)

Uses *cat* as mapper and *wc* as reducer:

hadoop jar \
\$HADOOP_HOME/share/hadoop/tools/lib/hadoopstreaming*.jar \
-input /datasets/plays/ -output streaming-out \
-mapper '/bin/cat' -reducer '/usr/bin/wc -l

Hadoop Ecosystem Lives On











Spark Capabilities

(i.e. Hadoop shortcomings)

- Performance
 - First, use RAM
 - Also, be smarter
- Ease of Use
 - Python, Scala, Java first class citizens
- New Paradigms
 - SparkSQL
 - Streaming
 - MLib
 - GraphX
 - ...more

But using Hadoop as the backing store is a common and sensible option.

Same Idea (improved)



RDD Resilient Distributed Dataset

Spark Formula

1. Create/Load RDD

Webpage visitor IP address log

2. Transform RDD

"Filter out all non-U.S. IPs"

3. But don't do anything yet!

Wait until data is actually needed Maybe apply more transforms ("Distinct IPs)

4. Perform Actions that return data

Count "How many unique U.S. visitors?"

Simple Example

>>> lines_rdd = sc.textFile("nasa_serverlog_20190404.tsv")



Spark Context

The first thing a Spark program requires is a context, which interfaces with some kind of cluster to use. Our pyspark shell provides us with a convenient *sc*, using the local filesystem, to start. Your standalone programs will have to specify one:

```
from pyspark import SparkConf, SparkContext
conf = SparkConf().setMaster("local").setAppName("Test_App")
sc = SparkContext(conf = conf)
```

```
You would typically run these scripts like so:
```

spark-submit Test_App.py

Simple Example

```
>>> lines_rdd = sc.textFile("nasa_serverlog_20190404.tsv")
>>> HubbleLines_rdd = lines_rdd.filter(lambda line: "Hubble" in line)
>>> HubbleLines_rdd.count()
47
>>> HubbleLines_rdd.first()
```



<u>Lambdas</u>

We'll see a lot of these. A lambda is simply a function that is too simple to deserve its own subroutine. Anywhere we have a lambda we could also just name a real subroutine that could go off and do anything.

When all you want to do is see if "given an input variable line, is "stanford" in there?", it isn't worth the digression.

Most modern languages have adopted this nicety.

'www.nasa.gov\shuttle/missions/61-c/Hubble.gif'

Common Transformations

Transformation	Result	
map(func)	Return a new RDD by passing each element through <i>func</i> .	Same Size
filter(func)	Return a new RDD by selecting the elements for which <i>func</i> returns true.	Fewer Elements
flatMap(func)	<i>func</i> can return multiple items, and generate a sequence, allowing us to "flatten" nested entries (JSON) into a list.	More Elements
distinct()	Return an RDD with only distinct entries.	
sample()	Various options to create a subset of the RDD.	
union(RDD)	Return a union of the RDDs.	
intersection(RDD)	Return an intersection of the RDDs.	
subtract(RDD)	Remove argument RDD from other.	
cartesian(RDD)	Cartesian product of the RDDs.	
parallelize(list)	Create an RDD from this (Python) list (using a spark context).	

Full list at http://spark.apache.org/docs/latest/api/python/reference/api/pyspark.RDD.html#pyspark.RDD

Common Actions

Action	Result
collect()	Return all the elements from the RDD.
count()	Number of elements in RDD.
countByValue()	List of times each value occurs in the RDD.
reduce(func)	Aggregate the elements of the RDD by providing a function which combines any two into one (sum, min, max,).
first(), take(n)	Return the first, or first n elements.
top(n)	Return the n highest valued elements of the RDDs.
takeSample()	Various options to return a subset of the RDD
saveAsTextFile(path)	Write the elements as a text file.
foreach(func)	Run the <i>func</i> on each element. Used for side-effects (updating accumulator variables) or interacting with external systems.

Full list at http://spark.apache.org/docs/latest/api/python/reference/api/pyspark.RDD.html#pyspark.RDD

Transformations vs. Actions

Transformations go from one RDD to another¹.

Actions bring some data back from the RDD.

Transformations are where the Spark machinery can do its magic with lazy evaluation and clever algorithms to minimize communication and parallelize the processing. You want to keep your data in the RDDs as much as possible.

Actions are mostly used either at the end of the analysis when the data has been distilled down (*collect*), or along the way to "peek" at the process (*count, take*).

¹ Yes, some of them also create an RDD (parallelize), but you get the idea.

Pair RDDs

• Key/Value organization is a simple, but often very efficient schema, as we mentioned in our NoSQL discussion.

• Spark provides special operations on RDDs that contain key/value pairs. They are similar to the general ones that we have seen.

 On the language (Python, Scala, Java) side key/values are simply tuples. If you have an RDD <u>all</u> of whose elements happen to be tuples of two items, it is a Pair RDD and you can use the key/value operations that follow.

Pair RDD Transformations

Transformation	Result
reduceByKey(func)	Reduce values using <i>func</i> , but on a key by key basis. That is, combine values with the same key.
groupByKey()	Combine values with same key. Each key ends up with a list.
sortByKey()	Return an RDD sorted by key.
mapValues(func)	Use <i>func</i> to change values, but not key.
keys()	Return an RDD of only keys.
values()	Return an RDD of only values.

Note that all of the regular transformations are available as well.
Pair RDD Actions

As with transformations, all of the regular actions are available to Pair RDDs, and there are some additional ones that can take advantage of key/value structure.

Action	Result
countByKey()	Count the number of elements for each key.
lookup(key)	Return all the values for this key.

Two Pair RDD Transformations

Transformation	Result
subtractByKey(otherRDD)	Remove elements with a key present in other RDD.
join(otherRDD)	Inner join: Return an RDD containing all pairs of elements with matching keys in self and other. Each pair of elements will be returned as a (k, (v1, v2)) tuple, where (k, v1) is in self and (k, v2) is in other.
leftOuterJoin(otherRDD)	For each element (k, v) in self, the resulting RDD will either contain all pairs (k, (v, w)) for w in other, or the pair (k, (v, None)) if no elements in other have key k.
rightOuterJoin(otherRDD)	For each element (k, w) in other, the resulting RDD will either contain all pairs (k, (v, w)) for v in this, or the pair (k, (None, w)) if no elements in self have key k.
cogroup(otherRDD)	Group data from both RDDs by key.

Joins Are Quite Useful

Any database designer can tell you how common joins are. Let's look at a simple example. We have (here we create it) an RDD of our top purchasing customers.

And an RDD with all of our customers' addresses.

To create a mailing list of special coupons for those favored customers we can use a join on the two datasets.

```
>>> best_customers_rdd = sc.parallelize([("Joe", "$103"), ("Alice", "$2000"), ("Bob", "$1200")])
```

```
>>> customer_addresses_rdd = sc.parallelize([("Joe", "23 State St."), ("Frank", "555 Timer Lane"), ("Sally", "44
Forest Rd."), ("Alice", "3 Elm Road"), ("Bob", "88 West Oak")])
```

```
>>> promotion_mail_rdd = best_customers_rdd.join(customer_addresses_rdd)
```

```
>>> promotion_mail_rdd.collect()
[('Bob', ('$1200', '88 west Oak')), ('Joe', ('$103', '23 State St.')), ('Alice', ('$2000', '3 Elm Road'))]
```

Shakespeare, a Data Analytics Favorite

Applying data analytics to the works of Shakespeare has become all the rage. Whether determining the legitimacy of his authorship (it wasn't Marlowe) or if Othello is actually a comedy (perhaps), or which word makes Macbeth so creepy ("the", yes) it is amazing how much publishable research has sprung from the recent analysis of 400 year old text.



We're going to do some exercises here using a text file containing all of his works.

Who needs this Spark stuff?

As we do our first Spark exercises, you might think of several ways to accomplish these tasks that you already know. For example, Python *Pandas* is a fine way to do our following problem, and it will probably work on your laptop reasonably well. But they do not scale well*.

However we are learning how to leverage scalable techniques that work on very big data. Shortly, we will encounter problems that are considerable in size, and you will leave this workshop knowing how to harness very large resources.

Searching the *Complete Works of William Shakespeare* for patterns is a lot different from searching the entire Web (perhaps as the 800TB *Common Crawl* dataset).

So everywhere you see an RDD, realize that it is a actually a parallel databank that could scale to PBs.



* See Panda's creator Wes McKinney's "10 Things I Hate About Pandas" at https://wesmckinney.com/blog/apache-arrow-pandas-internals/





Hands-On Time

We have an input file, Complete _Shakespeare.txt, that you can also find at <u>http://www.gutenberg.org/ebooks/100</u>.

You might find it useful to have <u>http://spark.apache.org/docs/latest/api/python/reference/api/pyspark.RDD.html#pyspark.RDD</u> in a browser window.

If you are starting from scratch on the login node, the first time you should grab all of our big data exercise sets and copy them to your home directory:

- 1) cp -r ~training/BigData
- 2) interact
- 3) cd BigData/Shakespeare
- 4) module load spark

```
5) pyspark
...
rdd = sc.textFile("Complete_Shakespeare.txt")
rdd.count()
Out: 124787
```

This last count command is a useful check because Spark is pretty quite about errors which, combined with its lazy evaluation, means you could misspell the filename, or be in the wrong directory and you won't find out until a later command.

Some Simple Problems (20 minutes?)

Let's try a few simple exercises.

- 1) Count the number of lines
- Count the number of words (hint: Python "split" is a workhorse)
- 3) Count unique words
- 4) Count the occurrence of each word
- 5) Show the top 5 most frequent words

These last two are a bit more challenging. One approach is to think "key/value". If you go that way, think about which data should be the key and don't be afraid to swap it about with value. This is a very common manipulation when dealing with key/value organized data.

Some Simple Answers

```
>>> lines_rdd = sc.textFile("Complete_Shakespeare.txt")
>>> lines_rdd.count()
124787
>>> words_rdd = lines_rdd.flatMap(lambda x: x.split())
>>> words_rdd.count()
904061
>>>
>>> words_rdd.distinct().count()
67779
>>>
```

Next, I know I'd like to end up with a pair RDD of sorted word/count pairs:

```
(23407, 'the'), (19540,'I'), (15682, 'to'), (15649, 'of') ...
```

Why not just *words_rdd.countByValue()*? It is an *action* that gives us a massive Python unsorted dictionary of results:

... 1, 'precious-princely': 1, 'christenings?': 1, 'empire': 11, 'vaunts': 2, 'Lubber's': 1, 'poet.': 2, 'Toad!': 1, 'leaden': 15, 'captains': 1, 'leaf': 9, 'Barnes,': 1, 'lead': 101, 'Hell': 1, 'wheat,': 3, 'lean': 28, 'Toad,': 1, 'trencher!': 2, '1.F.2.': 1, 'leas': 2, 'leap': 17, ...

Where to go next? Sort this in Python or try to get back into an RDD? If this is truly *BIG* data, we want to remain as an RDD until we reach our final results. So, no.

Some Harder Answers



results_rdd = lines_rdd.flatMap(lambda x: x.split()).map(lambda x: (x,1)).reduceByKey(lambda x,y: x+y).map(lambda x: (x[1],x[0])).sortByKey(False)

Spark Anti-Patterns

Here are a couple code clues that you are not working with Spark, but probably against it.

```
for loops, collect in middle of analysis, large data structures
....
intermediate_results = data_rdd.collect()
python_data = []
for datapoint in intermediate_results:
    python_data.append(modify_datapoint(datapoint))
next_rdd = sc.parallelize(python_data)
....
```

Ask yourself, "would this work with billions of elements?". And likely anything you are doing with a for is something that Spark will gladly parallelize for you, if you let it.

Some Homework Problems

To do research-level text analysis, we generally want to clean up our input. Here are some of the kinds of things you could do to get a more meaningful distinct word count.

1) Remove punctuation. Often punctuation is just noise, and it is here. Do a Map and/or Filter (some punctuation is attached to words, and some is not) to eliminate all punctuation from our Shakespeare data. Note that if you are familiar with regular expressions, Python has a ready method to use those.

2) Remove stop words. Stop words are common words that are also often uninteresting ("I", "the", "a"). You can remove many obvious stop words with a list of your own, and the *MLlib* that we are about to investigate has a convenient *StopWordsRemover()* method with default lists for various languages.

3) Stemming. Recognizing that various different words share the same root ("run", "running") is important, but not so easy to do simply. Once again, Spark brings powerful libraries into the mix to help. A popular one is the Natural Language Tool Kit. You should look at the docs, but you can give it a quick test quite easily:

```
import nltk
from nltk.stem.porter import *
stemmer = PorterStemmer()
stems_rdd = words_rdd.map( lambda x: stemmer.stem(x) )
```

Optimizations

We said one of the advantages of Spark is that we can control things for better performance. There are a multitude of optimization, performance, tuning and programmatic features to enable better control. We quickly look at a few of the most important.

- Persistence
- Partitioning
- Parallel Programming Capabilities
- Performance and Debugging Tools

Persistence

- Lazy evaluation implies by default that all the RDD dependencies will be computed when we call an action on that RDD.
- If we intend to use that data multiple times (say we are filtering some log, then dumping the results, but we will analyze it further) we can tell Spark to persist the data.
- We can specify different levels of persistence: *MEMORY_ONLY, MEMORY_ONLY_SER, MEMORY_AND_DISK, MEMORY_AND_DISK_SER, DISK_ONLY*

Partitions

- Spark distributes the data of your RDDs across its resources. It tries to do some obvious things.
- With key/value pairs we can help keep that data grouped efficiently.
- We can create custom partitioners that beat the default (which is probably a hash or maybe range).
- Use persist() if you have partitioned your data in some smart way. Otherwise it will keep getting re-partitioned.

Parallel Programming Features

Spark has several parallel programming features that make it easier and more efficient to do operations in parallel in a more explicit way.

Accumulators are variables that allow many copies of a variable to exist on the separate worker nodes.

It is also possible to have replicated data that we would like all the workers to have access to. Perhaps a lookup table of IP addresses to country codes so that each worker can transform or filter on such information. Maybe we want to exclude all non-US IP entries in our logs. You might think of ways you could do this just by passing variables, but they would likely be expensive in actual operation (usually requiring multiple sends). The solution in Spark is to send an (immutable, read only) broadcast variable

Accumulators log = sc.textFile("logs") blanks = sc.accumlator(0) def tokenizeLog(line) global blanks # write-only variable if (line =="") blanks += 1 return line.split(" ") entries = log.flatMap(tokenizeLog)

entries.saveAsTextFile("parsedlogs.txt") print "Blank entries: %d" blanks.value

Broadcast Variables

```
log = sc.textFile("log.txt")
```

```
IPtable = sc.broadcast(loadIPTable())
```

Performance & Debugging

We will give unfortunately short shrift to performance and debugging, which are both important. Mostly, this is because they are very configuration and application dependent.

Here are a few things to at least be aware of:

- SparkConf() class. A lot of options can be tweaked here.
- Spark Web UI. A very friendly way to explore all of these issues.

IO Formats

Spark has an impressive, and growing, list of input/output formats it supports. Some important ones:

- Text
- CSV
- SQL type Query/Load
 - JSON (can infer schema)
 - Parquet
 - Hive
 - XML
 - Sequence (Hadoopy key/value)
 - Databases: JDBC, Cassandra, HBase, MongoDB, etc.
- Compression (gzip...)

And it can interface directly with a variety of filesystems: local, HDFS, Lustre, Amazon S3,...

Spark Streaming

Spark addresses the need for streaming processing of data with a API that divides the data into batches, which are then processed as RDDs.

There are features to enable:

- (quantification of the amount of data Fast recovery from 0
- Load balancing 0
- 0
- 0

15% of the "global datasphere" created, captured, and replicated across Integration with sta the world) is currently real-time. That Integration with oth number is growing quickly both in absolute terms and as a percentage.

A Few Words About DataFrames

As mentioned earlier, an appreciation for having some defined structure to your data has come back into vogue. For one, because it simply makes sense and naturally emerges in many applications. Often even more important, it can greatly aid optimization, especially with the Java VM that Spark uses.

For both of these reasons, you will see that the newest set of APIs to Spark are DataFrame based. Sound leading-edge? This is simply SQL type columns. Very similar to Python pandas DataFrames (but based on RDDs, so not exactly).

We haven't prioritized them here because they aren't necessary, and some of the pieces aren't mature. But some of the latest features use them.

And while they would just complicate our basic examples, they are often simpler for real research problems. So don't shy away from using them.

Creating DataFrames

It is very pretty intuitive to utilize DataFrames. Your elements just have labeled columns.

A *row RDD* is the basic way to go from RDD to DataFrame, and back, if necessary. A "row" is just a tuple.

Creating DataFrames

You will come across DataFrames created without a schema. They get default column names.

```
>>> noSchemaDataFrame = spark.createDataFrame( row_rdd )
>>> noSchemaDataFrame.show()
+----+---+--+
| _1| _2| _3| _4|
+----+---+
| Joe|Pine St.| PA|12543|
|Sally| Fir Dr.| WA|78456|
| Jose| Elm Pl.| ND|45698|
+----+---++
```

Datasets

Spark has added a variation (technically a superset) of *DataFrames* called *Datasets*. For compiled languages with strong typing (Java and Scala) these provide static typing and can detect some errors at compile time.

This is not relevant to Python or R.

And you can create them inline as well.

Just Spark DataFrames making life easier...

Data from https://github.com/spark-examples/pyspark-examples/raw/master/resources/zipcodes.json

{"RecordNumber":1,"ZipCode":704,"ZipCodeType":"STANDARD","City":"PARC PARQUE","State":"PR","LocationType":"NOT ACCEPTABLE","Lat":17.96,"Long":-66.22,"Xaxis":0.38,"Yaxis":-0.87,"Zaxis":0.3,"WorldRegion": {"RecordNumber":2,"ZipCode":704,"ZipCodeType":"STANDARD","City":"PASEO COSTA DEL SUR","State":"PR","LocationType":"NOT ACCEPTABLE","Lat":17.96,"Long":-66.22,"Xaxis":0.38,"Yaxis":-0.87,"Zaxis":0.3,"WorldRegion {"RecordNumber":10,"ZipCode":709,"ZipCodeType":"STANDARD","City":"BDA SAN LUIS","State":"PR","LocationType":"NOT ACCEPTABLE","Lat":18.14,"Long":-66.26,"Xaxis":0.38,"Yaxis":-0.86,"Zaxis":0.31,"WorldRegion

<pre>>>> df = spark.read.json("zipcodes.json")</pre>	<pre>>>> df.show()</pre>					
>>> df.printSchema()	+	+	+	+	+	+
root	City	Country	Decommisioned	EstimatedPopulation	Lat	Location
City: string (nullable = true)	+	+	+	+	+	+
Country: string (nullable = true)	PARC PARQUE	US	false	null	17.96	NA-US-PR-PARC PARQUE
<pre> Decommisioned: boolean (nullable = true)</pre>	PASEO COSTA DEL SUR	US	false	null	17.96	NA-US-PR-PASEO CO
<pre> EstimatedPopulation: long (nullable = true)</pre>	BDA SAN LUIS	US	false	null	18.14	NA-US-PR-BDA SAN
Lat: double (nullable = true)	CINGULAR WIRELESS	US	false	null	32.72	NA-US-TX-CINGULAR
<pre> Location: string (nullable = true)</pre>	FORT WORTH	US	false	4053	32.75	NA-US-TX-FORT WORTH
<pre> LocationText: string (nullable = true)</pre>	FT WORTH	US	false	4053	32.75	NA-US-TX-FT WORTH
<pre> LocationType: string (nullable = true)</pre>	URB EUGENE RICE	US	false	null	17.96	NA-US-PR-URB EUGE
Long: double (nullable = true)	MESA	US	false	26883	33.37	NA-US-AZ-MESA
<pre> Notes: string (nullable = true)</pre>	MESA	US	false	25446	33.38	NA-US-AZ-MESA
RecordNumber: long (nullable = true)	HILLIARD	US	false	7443	30.69	NA-US-FL-HILLIARD
State: string (nullable = true)	HOLDER	US	false	null	28.96	NA-US-FL-HOLDER
TaxReturnsFiled: long (nullable = true)	HOLT	US	false	2190	30.72	NA-US-FL-HOLT
TotalWages: long (nullable = true)	HOMOSASSA	US	false	null	28.78	NA-US-FL-HOMOSASSA
<pre> WorldRegion: string (nullable = true)</pre>	BDA SAN LUIS	US	false	null	18.14	NA-US-PR-BDA SAN
Xaxis: double (nullable = true)	SECT LANAUSSE	US	false	null	17.96	NA-US-PR-SECT LAN
Yaxis: double (nullable = true)	SPRING GARDEN	US	false	null	33.97	NA-US-AL-SPRING G
Zaxis: double (nullable = true)	SPRINGVILLE	US	false	7845	33.77	NA-US-AL-SPRINGVILLE
ZipCodeType: string (nullable = true)	SPRUCE PINE	US	false	1209	34.37	NA-US-AL-SPRUCE PINE
Zipcode: long (nullable = true)	ASH HILL	US	false	1666	36.4	NA-US-NC-ASH HILL
	ASHEBORO	US	false	15228	35.71	NA-US-NC-ASHEBORO

And Sometime DataFrames Are Limiting

DataFrames are not as flexible as plain RDDs, and it isn't uncommon to find yourself fighting to do something that would be simple with a map, for example. In that case, don't hesitate to flip back into a plain RDD.

```
>>> aDataFrameFromRDD = spark.createDataFrame( row_rdd, ["name", "street", "state", "zip"] )
```

```
>>> another_row_rdd = aDataFrameFromRDD.rdd
```

Notice that this is not even a method, it is just a property. This is a clue that behind the scenes we are always working with RDDs.

A minor technicality here is that the returned object is actually a "Row" type. You may not care. If you want it be the original tuple type then

```
>>> tuple_rdd = aDataFrameFromRDD.rdd.map(tuple)
```

Note that when our map function is a function that already exists, there is no need for a lambda.

Speaking of pandas, or SciPy, or...

Some of you may have experience with the many Python libraries that accomplish some of these tasks. Immediately relevant to today, *pandas* allows us to sort and query data, and *SciPy* provides some nice clustering algorithms. So why not just use them?

The answer is that Spark does these things in the context of having potentially huge, parallel resources at hand. We don't notice it as Spark is also convenient, but behind every Spark call:

- every RDD could be many TB in size
- every transform could use many thousands of cores and TB of memory
- every algorithm could also use those thousands of cores

So don't think of Spark as just a data analytics library because our exercises are modest. You are learning how to cope with Big Data.

Other Scalable Alternatives: Dask

Of the many alternatives to play with data on your laptop, there are only a few that aspire to scale up to big data. The only one, besides Spark, that seems to have any traction is Dask.

It attempts to retain more of the "laptop feel" of your toy codes, making for an easier port. The tradeoff is that the scalability is a lot more mysterious. If it doesn't work - or someone hasn't scaled the piece you need - your options are limited.

At this time, I'd say it is riskier, but academic projects can often entertain more risk than industry.

Numpy like operations

Dataframes implement Pandas

import dask.dataframe as dd
df = dd.read_csv('/.../2020-*-*.csv')
df.groupby(df.account_id).balance.sum()

Pieces of Scikit-Learn

from dask_ml.linear_model import \
LogisticRegression
lr = LogisticRegression()
lr.fit(train, test)

Drill Down?

Using MLlib

One of the reasons we use spark is for easy access to powerful data analysis tools. The MLlib library gives us a machine learning library that is easy to use and utilizes the scalability of the Spark system.

It has supported APIs for Python (with NumPy), R, Java and Scala.

We will use the Python version in a generic manner that looks very similar to any of the above implementations.

There are good example documents for the clustering routine we are using, as well as alternative clustering algorithms, here:

http://spark.apache.org/docs/latest/mllib-clustering.html

I suggest you use these pages for your Spark work.

Clustering

Clustering is a very common operation for finding grouping in data and has countless applications. This is a very simple example, but you will find yourself reaching for a clustering algorithm frequently in pursuing many diverse machine learning objectives, sometimes as one part of a pipeline.



Weight

Coin Sorting

Clustering

As intuitive as clustering is, it presents challenges to implement in an efficient and robust manner.



We will start with 5000 2D points. We want to figure out how many clusters there are, and their centers. Let's fire up pyspark and get to it...



Using Python version 2.7.5 (default, Nov 20 2015 02:00:19) SparkContext available as sc, HiveContext available as sqlContext. >>> >>> rdd1 = sc.textFil br06% interact to RDD >>> >>> rdd2 = rdd1.map(] rm to words and integers >>> rdd3 = rdd2.map(1)>>> r288% r288% module load spark r288% pyspark

Finding Our Way

```
>>> rdd1 = sc.textFile("5000_points.txt")
>>> rdd1.count()
5000
>>> rdd1.take(4)
Γ'
      664159
               550946', '
                           665845
                                       557965', 597173
                                                               575538'. '
                                                                             618600
                                                                                      551446']
>>> rdd2 = rdd1.map(lambda x:x.split())
>>> rdd2.take(4)
[['664159', '550946'], ['665845', '557965'], ['597173', '575538'], ['618600', '551446']]
>>> rdd3 = rdd2.map(lambda x: [int(x[0]),int(x[1])])
>>> rdd3.take(4)
[[664159, 550946], [665845, 557965], [597173, 575538], [618600, 551446]]
>>>
```

Using Python version 2.7.5 (default, Nov 20 2015 02:00:19)
SparkContext available as sc, HiveContext available as sqlContext.
>>>
rdd1 = sc.textFile("5000_points.txt")
>>>
rdd2 = rdd1.map(lambda x:x.split())
>>> rdd3 = rdd2.map(lambda x: [int(x[0]),int(x[1])])
>>>
>>>

>>> from pyspark.mllib.clustering import KMeans







```
Using Python version 2.7.5 (default, Nov 20 2015 02:00:19)
SparkContext available as sc, HiveContext available as sqlContext.
>>>
>>> rdd1 = sc.textFile("5000_points.txt")
>>>
>>> rdd2 = rdd1.map(lambda x:x.split())
>>> rdd3 = rdd2.map(lambda x: [int(x[0]),int(x[1])])
>>>
    from pyspark.mllib.clustering import KMeans
>>>
>>>
>>> for clusters in range(1,30):
        model = KMeans.train(rdd3, clusters)
                                                            Let's see results for 1-30 cluster tries
        print (clusters, model.computeCost(rdd3))
. . .
. . .
```

1 5.76807041184e+14 2 3.43183673951e+14 3 2.23097486536e+14 4 1.64792608443e+14 5 1.19410028576e+14 6 7.97690150116e+13 7 7.16451594344e+13 8 4.81469246295e+13 9 4.23762700793e+13 10 3.65230706654e+13 11 3.16991867996e+13 12 2.94369408304e+13 13 2.04031903147e+13 14 1.37018893034e+13 15 8.91761561687e+12 16 1.31833652006e+13 17 1.39010717893e+13 18 8.22806178508e+12 19 8.22513516563e+12 20 7.79359299283e+12 21 7.79615059172e+12 22 7.70001662709e+12 23 7.24231610447e+12 24 7.21990743993e+12 25 7.09395133944e+12 26 6.92577789424e+12 27 6.53939015776e+12 28 6.57782690833e+12 29 6.37192522244e+12

Right Answer?

	12 2.45472346524e+13	12 2.31466520037e+13
	13 2.00175423869e+13	13 1.91856542103e+13
<pre>>>> for trials in range(10):</pre>	14 1.90313863726e+13	14 1.49332023312e+13
nrint	15 1.52746006962e+13	15 1.3506302755e+13
for all store in $range(12, 10)$.	16 8.67526114029e+12	16 8.7757678836e+12
For Clusters in range(12,18):	17 8.49571894386e+12	17 1.60075548613e+13
<pre>model = KMeans.train(rdd3,clusters)</pre>		
<pre>print (clusters, model.computeCost(rdd3))</pre>	12 2.62619056924e+13	12 2.5187054064e+13
	13 2.90031673822e+13	13 1.83498739266e+13
	14 1.52308079405e+13	14 1.96076943156e+13
	15 8.91765957989e+12	15 1.41725666214e+13
	16 8.70736515113e+12	16 1.41986217172e+13
	17 8.49616440477e+12	17 8.46755159547e+12
	12 2.5524719797e+13	12 2.38234539188e+13
	13 2.14332949698e+13	13 1.85101922046e+13
	$14 \ 2 \ 11070395905e+13$	14 1 91732620477e+13
	15 1.47792736325e+13	15 8.91769396968e+12
	16 1 85736955725e+13	16 8 64876051004e+12
	17 8.42795740134e+12	17 8.54677681587e+12
		2. 0.0.00000000000000000000000000000000
	12 2 31466242693e+13	$12 \ 2 \ 5187054064e+13$
	13 2.10129797745e+13	13 2.04031903147e+13
	14 1.45400177021e+13	14 1.95213876047e+13
	15 1.52115329071e+13	15 1.93000628589e+13
	$16 \ 1 \ 41347332901e+13$	16 2 07670831868e+13
	$17 \ 1 \ 31314086577e+13$	17 8 47797102908e+12
	1/ 1/5151/0005//0115	1/ 0/ // 5/ 1025000/12
	12 2 47927778784e+13	12 2 39830397362e+13
	13 2.43404436887e+13	13 2.00248378195e+13
	$14, 2, 1522702068e \pm 13$	14 1,34867337672e+13
	15 8 91765000665 <u>e+12</u>	15 2 09299321238e+13
	$16 \ 1 \ 4580927737e \pm 13$	16 1 32266735736e+13
	$17 8 57823507015_{P+12}$	17 8 50857884943e+12

Find the Centers

>>>	for trials in range(10):	<pre>#Try ten times to find best result</pre>
	for clusters in range(12, 16):	#Only look in interesting range
	<pre>model = KMeans.train(rdd3, clusters)</pre>	
	<pre>cost = model.computeCost(rdd3)</pre>	
	<pre>centers = model.clusterCenters</pre>	#Let's grab cluster centers
	if cost<1e+13:	#If result is good, print it out
	print (clusters, cost)	
	for coords in centers:	
	<pre>print (int(coords[0]), int(coords[1]))</pre>	
	break	

398555 404855 417799 787001



♦ Series1
16 Clusters



Dimensionality Reduction

We are going to find a recurring theme throughout machine learning:

- Our data naturally resides in higher dimensions
- Reducing the dimensionality makes the problem more tractable
- And simultaneously provides us with insight

This last two bullets highlight the principle that "learning" is often finding an effective compressed representation.

As we return to this theme, we will highlight these slides with our Dimensionality Reduction badge so that you can follow this thread and appreciate how fundamental it is.



Why all these dimensions?



The problems we are going to address, as well as the ones you are likely to encounter, are naturally highly dimensional. If you are new to this concept, lets look at an intuitive example to make it less abstract.

Category	Purchase Total (\$)	
Children's Clothing	\$800	^
Pet Supplies	\$0	
Cameras (Dash, Security, Baby)	\$450	
Containers (Storage)	\$350	UF
Romance Book	\$0	C
Remodeling Books	\$80	
Sporting Goods	\$25	Ite
Children's Toys	\$378	ju i i i i i i i i i i i i i i i i i i i
Power Tools	\$0	or
Computers	\$0	le c
Garden	\$0	- 01
Children's Books	\$180	V

Why all these dimensions?



If we apply our newfound clustering expertise, we might find we have 80 clusters (with an acceptable error).

People spending on "child's toys " and "children's clothing" might cluster with "child's books" and, less obvious, "cameras (Dashcams, baby monitors and security cams)", because they buy new cars and are safety conscious. We might label this cluster "Young Parents". We also might not feel obligated to label the clusters at all. We can now represent any customer by their distance from these 80 clusters.

					80 dimensional	vector			
Customer Representation									
Cluster	Young Parents	College Athlete	Auto Enthusiast	Knitter	Steelers Fan	Shakespeare Reader	Sci-Fi Fan	Plumber	
Distance	0.02	2.3	1.4	8.4	2.2	14.9	3.3	0.8	

We have now accomplished two things:

- we have compressed our data
- learned something about our customers (who to send a dashcam promo to).

Curse of Dimensionality



This is a good time to point out how our intuition can lead us astray as we increase the dimensionality of our problems - which we will certainly be doing - and to a great degree. There are several related aspects to this phenomenon, often referred to as the *Curse of Dimensionality*. One root cause of confusion is that our notion of Euclidian distance starts to fail in higher dimensions.



These plots show the distributions of pairwise distances between randomly distributed points within differently dimensioned unit hypercubes. Notice how all the points start to be about the same distance apart.

Once can imagine this makes life harder on a clustering algorithm!

There are other surprising effects: random vectors are almost all orthogonal; the unit sphere takes almost no volume in the unit square. These cause all kinds of problems when generalizing algorithms from our lowly 3D world.

Metrics



Even the definition of distance (the *metric*) can vary based upon application. If you are solving chess problems, you might find the Manhattan distance (or taxicab metric) to be most useful.



For comparing text strings, we might choose one of dozens of different metrics. For spell checking you might want one that is good for phonetic distance, or maybe edit distance. For natural language processing (NLP), you probably care more about tokens.

For genomics, you might care more about string sequences.

Some useful measures don't even qualify as metrics (usually because they fail the triangle inequality: $a + b \ge c$).

Alternative DR: Principal Component Analysis







3D Data Set

Maybe mostly 1D!

Alternative DR: Principal Component Analysis







Flatter 2D-ish Data Set

View down the 1st Princ. Comp.

Why So Many Alternatives?



Let's look at one more example today. Suppose we are tying to do a Zillow type of analysis and predict home values based upon available factors. We may have an entry (vector) for each home that captures this kind of data:

Home	• Data
Latitude	4833438 north
Longitude	630084 east
Last Sale Price	\$ 480,000
Last Sale Year	1998
Width	62
Depth	40
Floors	3
Bedrooms	3
Bathrooms	2
Garage	2
Yard Width	84
Yard Depth	60

There may be some opportunities to reduce the dimension of the vector here. Perhaps clustering on the geographical coordinates...

Principal Component Analysis Fail







1st Component Off Data Not Very Linear

D x W Is Not Linear But (DxW) Fits Well

Why the fascination with linear techniques?



The Streetlight Effect

This is a very real and powerful force throughout the sciences.

It is not because practitioners are dumb.

But, it is also very often neither explained nor justified.

Which leads to great confusion.

Why Would An Image Have 784 Dimensions?





greyscale images

0-																										0
																										o
																										0
								0	116	125	171	255	255	150	93											0
							0	169	253	253	253	253	253	253	218											o
5 -						0	169	253	253	253	213	142	176	253	253	122										0
					0	52	250	253	210	32	12	0	6	206	253	140										0
					0	77	251	210	25			o	122	248	253	65										o
						0	31	18				0	209	253	253	65										0
												117	247	253	198											0
10 -											76	247	253	231	63											o
										o	128	253	253	144	0											o
									0	176	246	253	159	12												0
									25	234	253	233	35													o
								0	198	253	253	141														0
15 -							0	78	248	253	189	12														0
							19	200	253	253	141															0
						0	134	253	253	173	12															o
						0	248	253	253	25																0
						0	248	253	253											37	150	150	150	147		0
20 -						0	248	253	253	253	253	253	253	253	168	143	166	253	253	253	253	253	253	253	123	0
						0	174	253	253	253	253	253	253	253	253	253	253	253	249	247	247	169	117	117	57	0
						0	0	118	123	123	123	166	253	253	253	155	123	123	41	0	0	0	0	0		0
								0	0	0	0	0	0	0	0	0	0	0								0
																										o
25 -																										0
																										0
																										0
	ò			5					10					15					20					25		

Central Hypothesis of Modern DL





Data Lives On A Lower Dimensional Manifold



Maybe Very Contiguous



Maybe A Small Set Of Disconnected

Images from Wikipedia

Testing These Ideas With Scikit-learn



import numpy as np import matplotlib.pyplot as plt from sklearn import (datasets, decomposition, manifold, random_projection) def draw(X, title): plt.figure() plt.xlim(X.min(0)[0],X.max(0)[0]); plt.ylim(X.min(0)[1],X.max(0)[1]) plt.xticks([]); plt.yticks([]) plt.title(title) for i in range(X.shape[0]): plt.text(X[i, 0], X[i, 1], str(y[i]), color=plt.cm.set1(y[i] / 10.))

digits = datasets.load_digits(n_class=6)
X = digits.data
y = digits.target

rp = random_projection.SparseRandomProjection(n_components=2, random_state=42)
X_projected = rp.fit_transform(X)
draw(X_projected, "Sparse Random Projection of the digits")

X_pca = decomposition.PCA(n_components=2).fit_transform(X)
draw(X_pca, "PCA (Two Components)")

tsne = manifold.TSNE(n_components=2, init='pca', random_state=0)
X_tsne = tsne.fit_transform(X)
draw(X_tsne, "t-SNE Embedding")

plt.show()

Sample of 64-dimensional digits dataset
01234501134501234505
44150522001321431314
34405345642225346001
0413510022101233344
45052200432434344344
50425450423450555041
35400222042333344450 52200432443431434465
31544222554403011345
51001120113333441505
12001311431314314053
13450123450555041354
00112011333344150512
44121554400123450123









How does all this fit together?

The Journey Ahead



As the Data Scientist wanders across the ill-defined boundary between Data Science and Machine Learning, in search of the fabled land of Artificial Intelligence, they find that the language changes from programming to a creole of linear algebra and probablity and statistics.