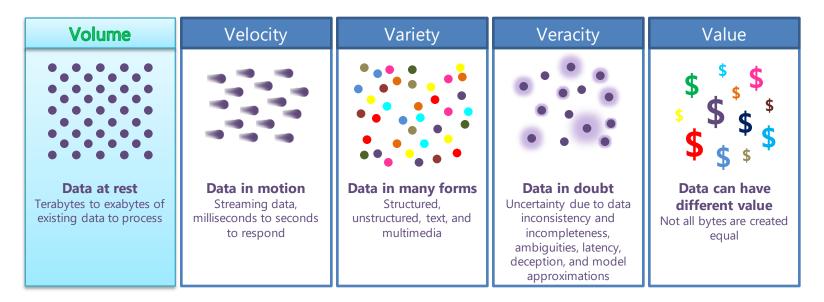
Big Data Engineering with Distributed Systems

Data Science Dojo



Batch Processing



- Save up all your raw data and process all at once
- Addresses the 'volume' problem of big data



Machine Learning Scaling

Programs	Programming	Cloud	Distributed	
• Excel	• Python • R	• Azure ML • AWS ML	HadoopSpark	
	• SAS	Big MLCloud Virtual Machines	• H20 • Revolution R	



Excel: Cell Meta Data

	Α	В	C	D	Е
1	Sepal.Leng	Sepal.Widt	Petal.Leng	Petal.Widt	Species
2	5.1	3.5	1.4	0.2	setosa
3	4.9	3	14	U.2	setosa

E2 Cell = Application, Address, AllowEdit, Areas, Borders, BottomPadding, Comment, Column, ColumnIndex, Creator Font, FitText, Height, HeightRule, ID, Interior, LeftPadding, NestingLevel, RightPadding, Row, RowIndex, Shading, Tables, TopPadding, VerticalAlignment, Value, Width, WordWrap "Font":{ "Application": "Microsoft Excel", "Background": None, "Bold": True, "Color": 0, "ColorIndex": 5. "Creator": "XCEL", "FontStyle": "Bold Italic", "Italic": True, "Name": "Comic Sans MS", "OutlineFont": True, "Parent": None, "Shadow": False, "Size": 12, "Strikethrough": False, "Subscript": False, "ThemeColor": 12, "ThemeFont:": 2, "TintAndShade": 1, "Superscript": False, "Underline": False,

"Value": "Setosa"



R Limits

- Single core
- Single threaded
- All in memory (RAM)
- Vectors & Matrices capped at 4,294,967,295
 elements (rows) if 32-bit version; 2^32 1

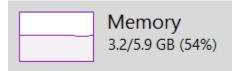


R Limits: RAM

All in memory (RAM)

 $Max\ Data\ Limit = (\ Total\ RAM\ Access\ x\ 80\%) - Normal\ RAM\ Usage$

Laptop Example:



 $Max\ Data\ Limit = (5.9\ GB\ x\ 80\%) - 3.2\ GB$ $Max\ Data\ Limit = \sim 1.52GB$

*R data frames actually bloats data files by 3x R's Data Limit = $1.52GB \div 3 = 518.8 MB$



R Limits: RAM

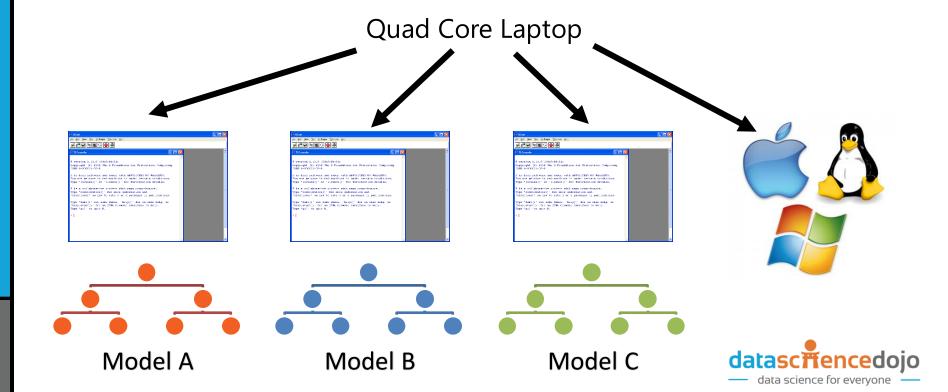
INSTANCE	CORES	RAM	DISK SIZES	PRICE
G1	2	28 GB	384 GB	\$0.67/hr (~\$498/mo)
G2	4	56 GB	768 GB	\$1.34/hr (~\$997/mo)
G3	8	112 GB	1,536 GB	\$2.68/hr (~\$1,994/mo)
G4	16	224 GB	3,072 GB	\$5.36/hr (~\$3,988/mo)
G5	32	448 GB	6,144 GB	\$9.65/hr (~\$7,180/mo)

Azure's Biggest Virtual Machine $Max\ Data\ Limit = (448gb\ x\ 80\%) - 1GB$ $Max\ Data\ Limit = \sim 357.4gb$



R Limits: Single Core

- Single core
- Single threaded



Machine Learning Scaling

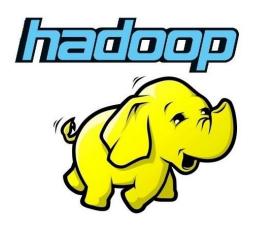
Programs	Programming	Cloud	Distributed
• Excel	PythonRSAS	Azure MLAWS MLWatson AnalyticsBig MLCloud Virtual Machines	 Hadoop Spark H20 Revolution R

Distributed R Solutions:

https://cran.r-project.org/web/views/HighPerformanceComputing.html



Agenda







From a Data Scientist's Perspective



Goals:

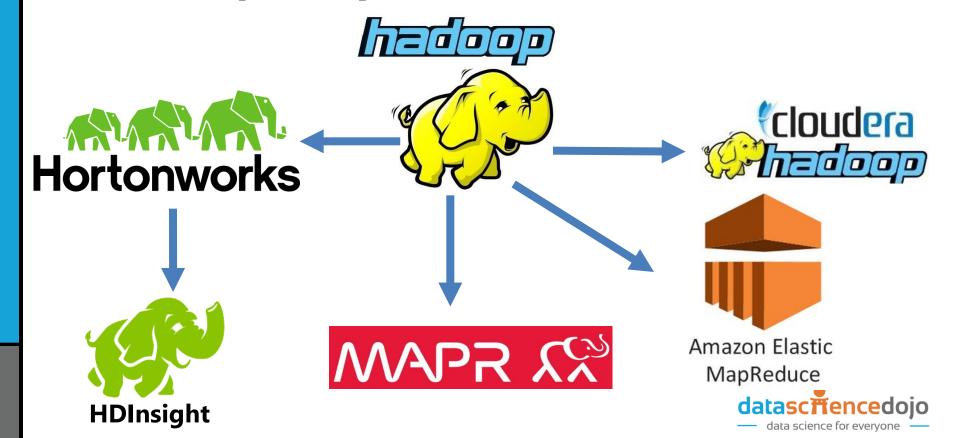
 Teach you how to leverage an existing Hadoop cluster, self-service data query

Not goals:

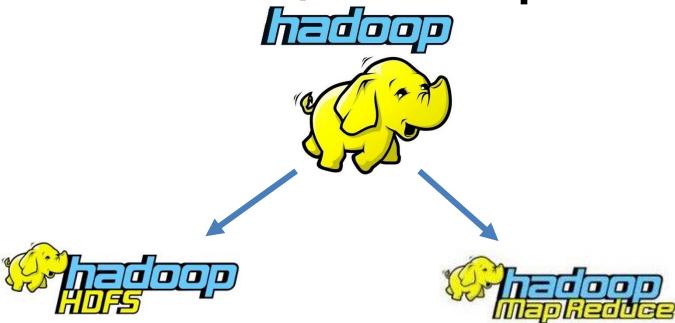
 Managing or administering a Hadoop cluster

data science for evervone

Hadoop Implementations



(Vanilla/Base) Hadoop



Processing engine for distributed batch processing.

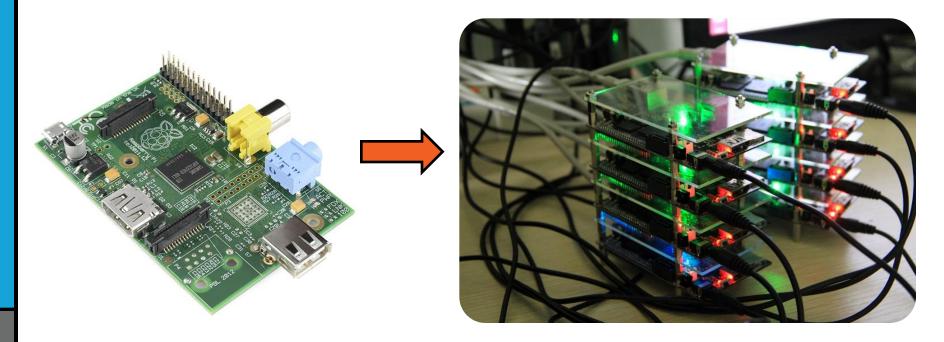


Turn Back The Clock, The Mainframe





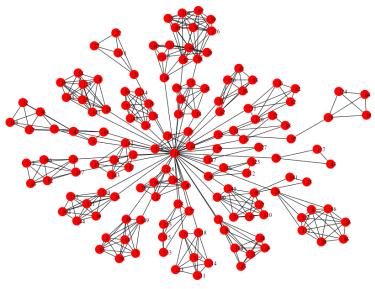
Distributed Computing





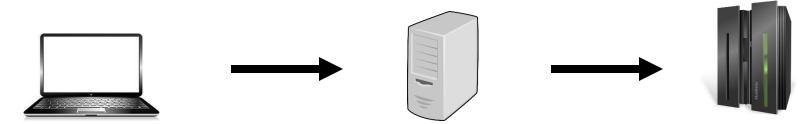
Cloud Computing







Scaling Computational Power



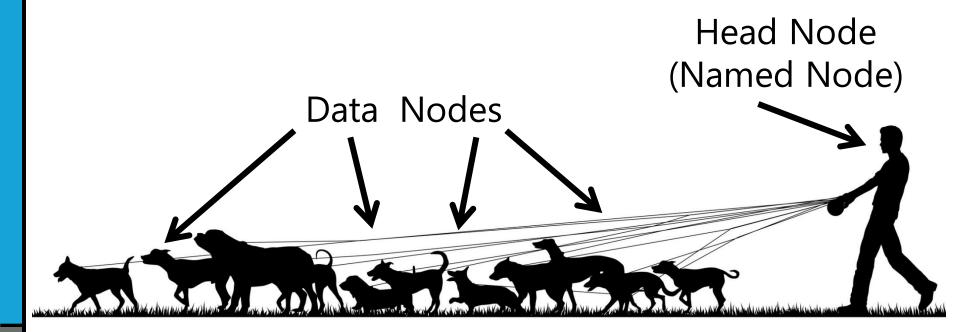
- Old Scaling:
- Vertical Scaling, Scaling UP
- High performance computers



- New Scaling:
- Horizontal Scaling, Scaling OUT
- Commodity hardware, distributed

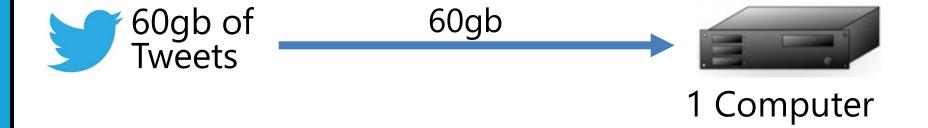


If dogs were servers...





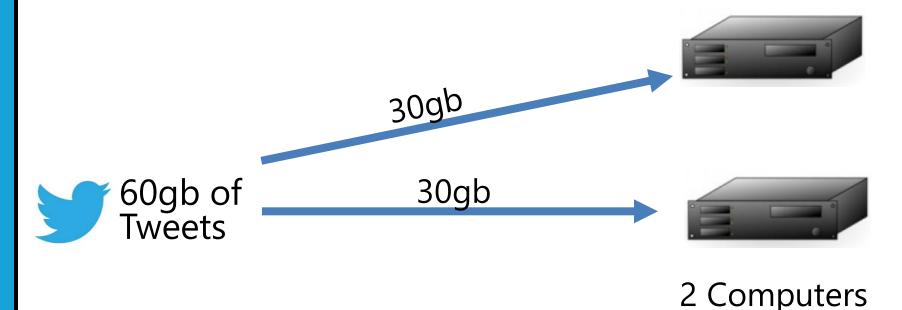
HDFS & MapReduce



Processing: 30 hours



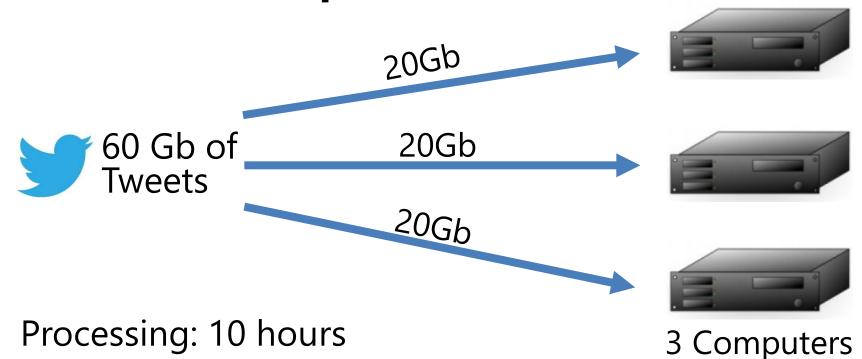
HDFS & MapReduce



Processing: 15 hours



HDFS & MapReduce





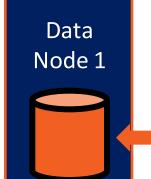
Most Cases, Linear Scaling Of Processing Power

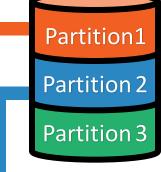
Number of Computers	Processing Time (hours)
1	30
2	15
3	10
4	7.5
5	6
6	5
7	4.26
8	3.75
9	3.33

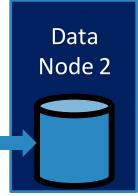


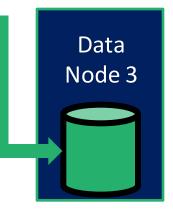
HDFS

HDFS Partitioning



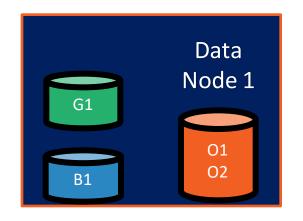


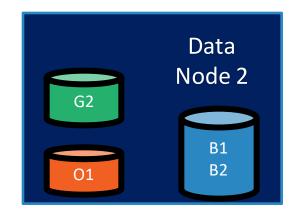


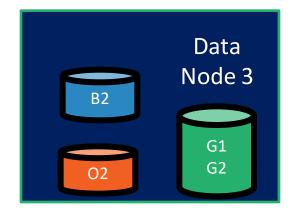




HDFS Redundancy









Limitations with MapReduce

- ~70 lines of code to do anything
- Slow
- Troubleshooting multiple computers
- Good devs are scarce
- Expensive certifications

```
org.apache.hadoop.examples;
import java.io.IOException;
import java.util.StringTokenizer;
       org.apache.hadoop.conf.Configuration;
       org.apache.hadoop.fs.Path;
       org.apache.hadoop.io.IntWritable;
       org.apache.hadoop.io.Text;
       org.apache.hadoop.mapreduce.Job;
       org.apache.hadoop.mapreduce.Mapper;
       org.apache.hadoop.mapreduce.Reducer;
       org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
       org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;
public class WordCount {
  public static class TokenizerMapper
       extends Mapper Object, Text, Text, IntWritable>{
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();
    public void map(Object key, Text value, Context context
                    ) throws IOException, InterruptedException {
      StringTokenizer itr = new StringTokenizer(value.toString());
      while (itr.hasMoreTokens()) {
        word.set(itr.nextToken());
        context.write(word, one);
```



Ambari: Cluster provisioning, management, and monitoring



Avro (Microsoft .NET Library for Avro): Data serialization for the Microsoft .NET environment



HBase: Non-relational database for very large tables



HDFS: Hadoop Distributed File System



Hive: SQL-like querying



Mahout: Machine learning





MapReduce and YARN: Distributed processing and resource management



Oozie: Workflow management



Pig: Simpler scripting for MapReduce transformations



Sqoop: Data import and export



STORM Storm: Real-time processing of fast, large data streams

Zookeeper: Coordinates processes in distributed systems

Hive Jobs

HiveQL Statement



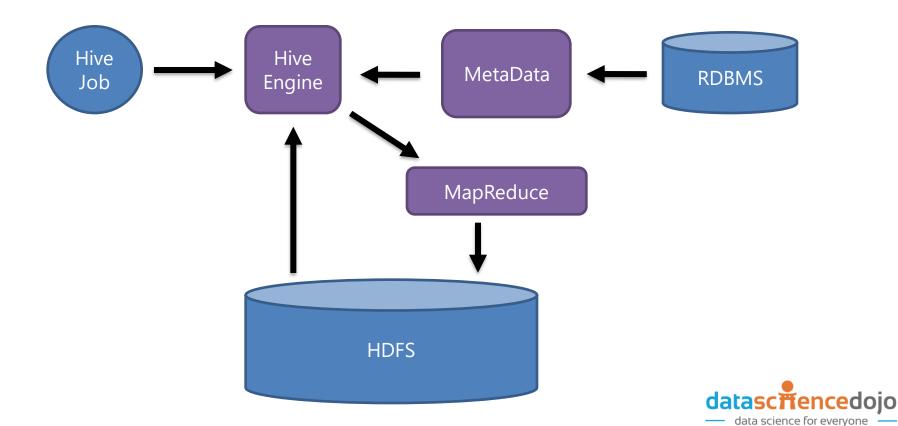
Translation & Conversion



MapReduce Job



Hive Architecture







Unstructured Data

Data File









Structured Data

Metadata File/DB



Semi Structured Data

Self Describing Flat Files

- XML
- JSON
- CSV
- TSV

```
"created_at": "Thu May 07 18:06:23 +0000 2015",
"id":596375540631646210,
"id_str": "596375540631646210",
"text": "Expert usable tips differently the pres:
"source": "<a href=\"http://twitterfeed.com\" rel
"truncated":0,
"in_reply_to_status_id":null,
"in_reply_to_status_id_str":null,
"in_reply_to_user_id":null,
"in_reply_to_user_id_str":null,
```



Why Hive?



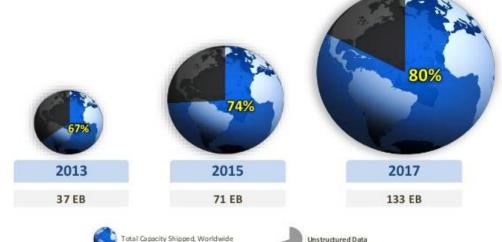
- SQL spoken here (HiveQL)
- ODBC driver
- BI Integration
- Supports only Structured Data



Limitations

Structured vs. Unstructured Data Growth





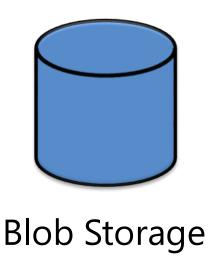






Azure Blob Storage







MapReduce, via Playing Cards

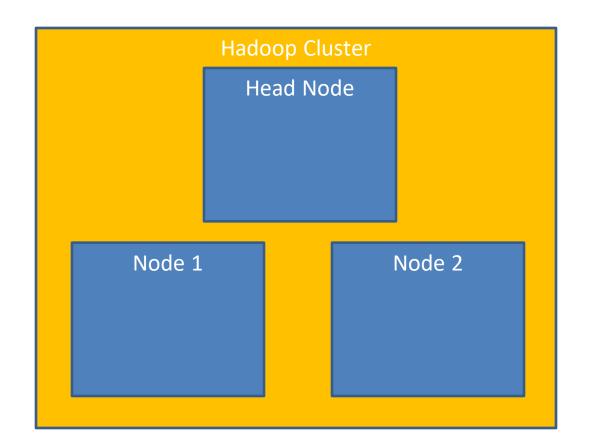


Let's count the number of spades, clubs, hearts, and diamonds in a stack of cards, the way map reduce would.

- Each card represents a row of data
- Each suit & number represents an attribute of the data

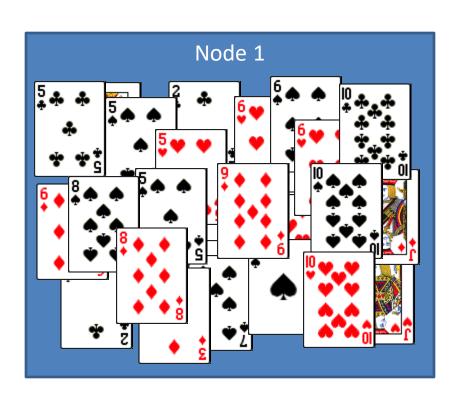


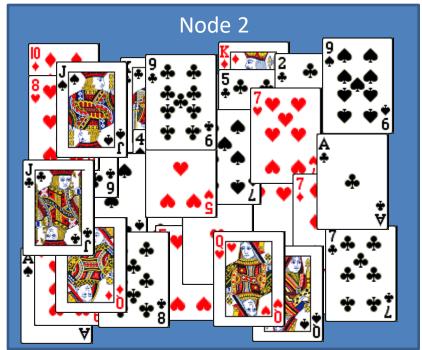
Using a 2 Data Node Cluster





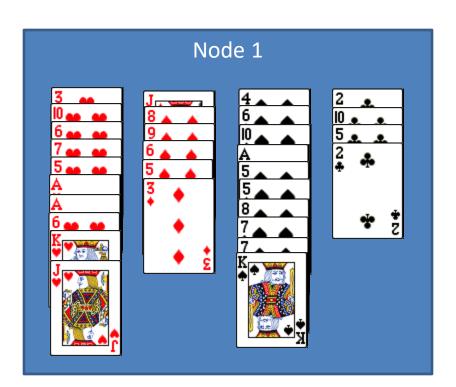
Mapping: Each Node's HDFS

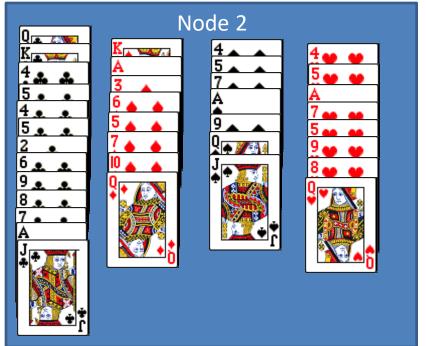






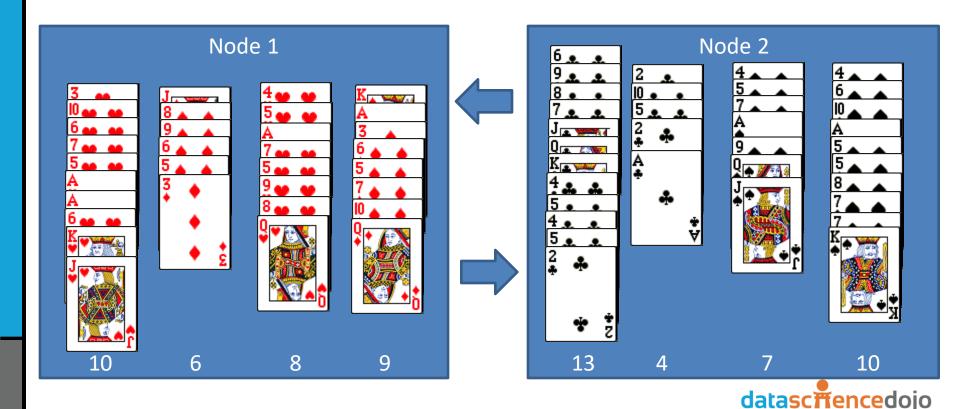
Mapping: Node Sorting





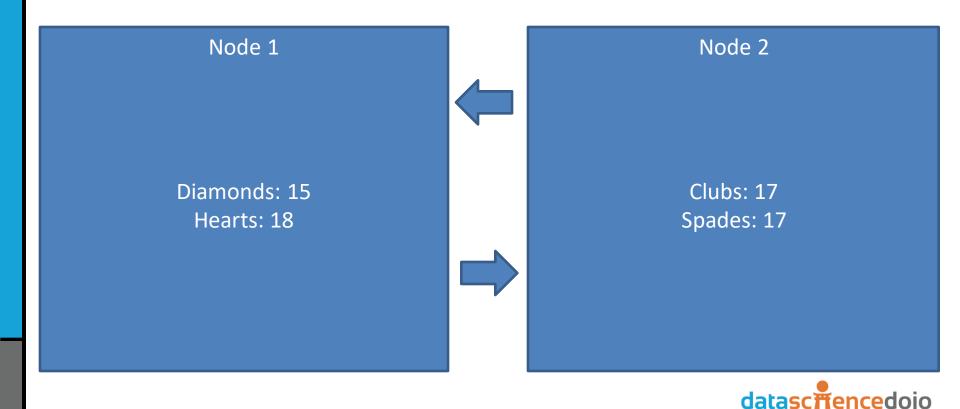


Mapping: Node Shuffle, Data Transfer



data science for everyone

Mapping: Node Shuffle, Data Transfer



data science for everyone

Execution Engine: Tez

The Stinger Initiative

2011, the world got together and declared MapReduce to be terrible.

- 44 companies
- 145 developers
- 392k lines of Java code

Hadoop 2.0 with Yarn & Tez

- Tez dropped hive query times by 90%, 100x performance
- Utilizes Apache Yarn
 - Yarn: resource manager for multi-cluster computing
- Introduced partial in-memory, local head nodes
- Rewrote HiveQL as an actual language, instead of translation

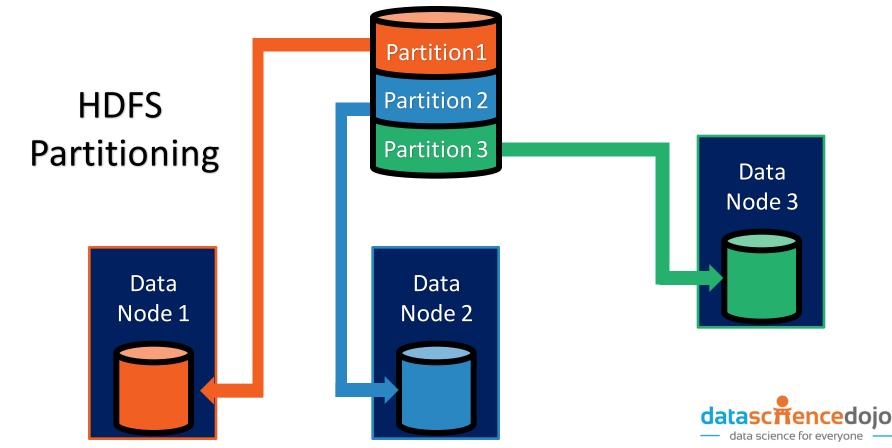




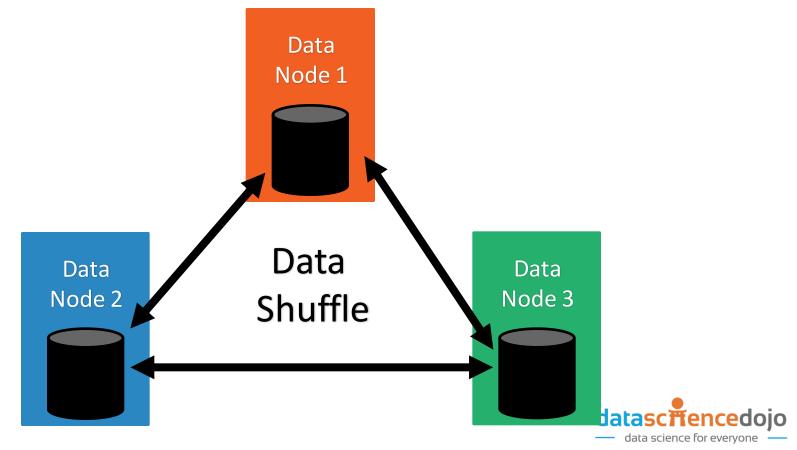
- Distributed Machine Learning
- Installed into Hadoop & Spark
- R-like language Implementation



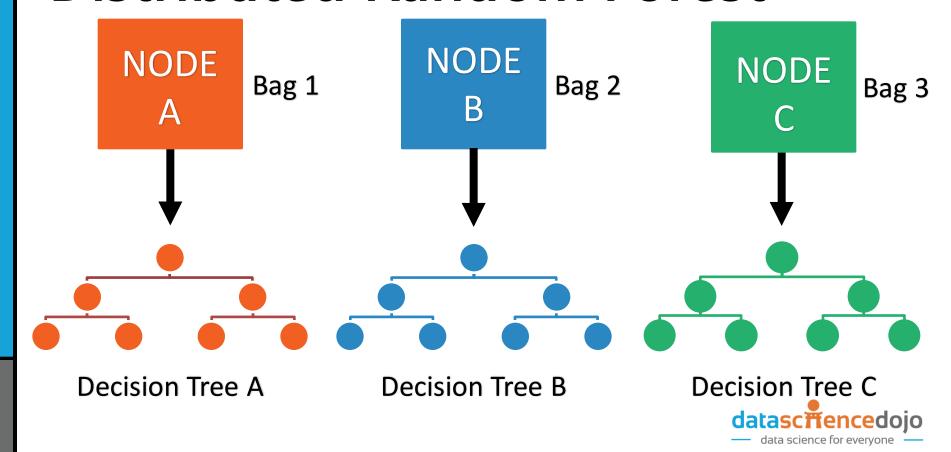
Distributed Random Forest



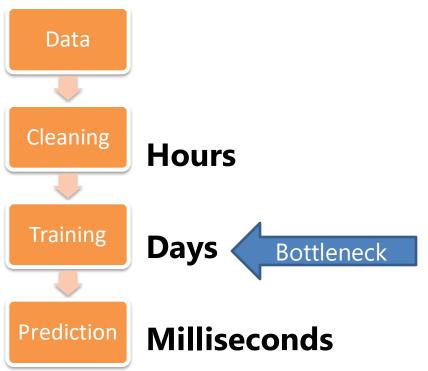
Distributed Random Forest



Distributed Random Forest



Processing Times - Machine Learning



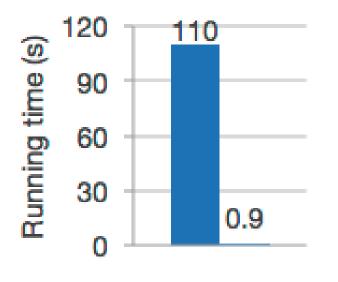
- Large scale systems are only needed for training
- Phones can use models outputted by mahout to predict new data
- After a model is trained, save the model to any IO file type and reload it where you want











Spark

In-Memory: 100x Hadoop times faster than Hadoop





3x faster on 10x few machines

Datona GraySort Benchmark: Sort 100 TB of data

Previous World Record: 2014:

- Method: Hadoop
- Yahoo!
- 72 Minutes
- 2100 Nodes

Method: Spark

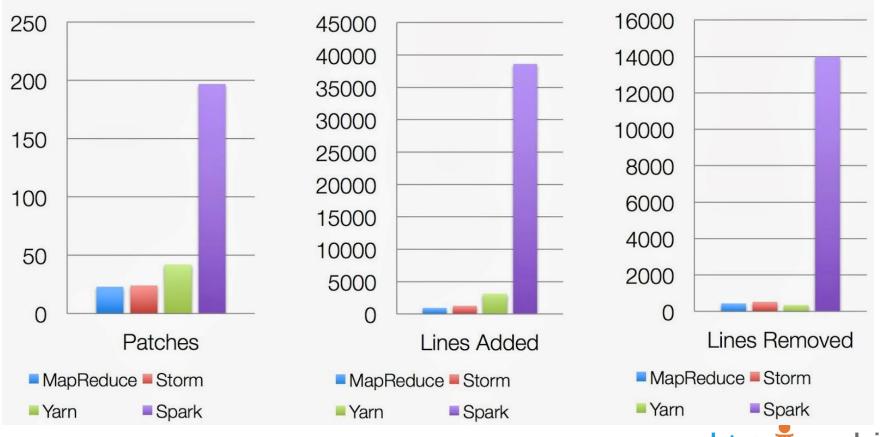
Databricks

23 Minutes

206 Nodes



Activity in last 30 days



Source: Xiangrui Meng, Data Bricks





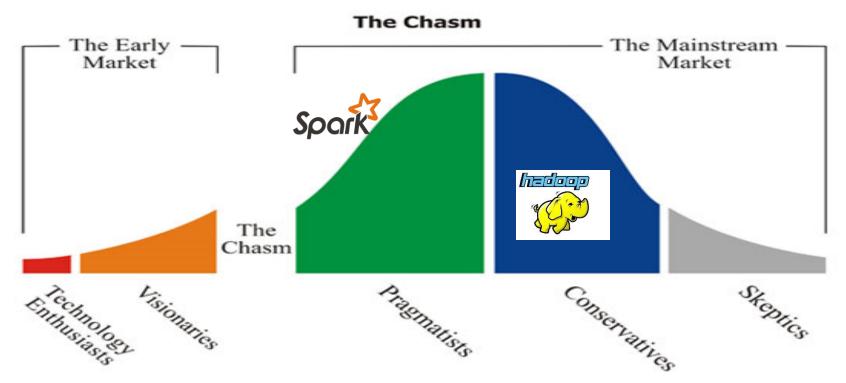
Spark SQL

Spark Streaming MLlib (machine learning) GraphX (graph)

Apache Spark



Technology adoption life cycle



Source: http://carlosmartinezt.com/2010/06/technology-adoption-life-cycle/



QUESTIONS



Enjoying the bootcamp?

We'd love it if you could write a short review of Data Science Dojo!

Switch Up (https://www.switchup.org/bootcamps/data-science-dojo)
Course Report (https://www.coursereport.com/schools/data-science-dojo)



datascmencedojo

data science for everyone

Your reviews help other people find and attend our bootcamp.