

Case Study 2: Document Retrieval

Parallel Programming Map-Reduce

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin

January 31st, 2013

©Carlos Guestrin 2013

1

Needless to Say, We Need Machine Learning for Big Data

flickr

6 Billion
Flickr Photos



28 Million
Wikipedia Pages

facebook

1 Billion
Facebook Users

You Tube

72 Hours a Minute
YouTube

The New York Times

SundayReview

WORLD U.S. N.Y. / REGION BUSINESS TEC

NEWS ANALYSIS

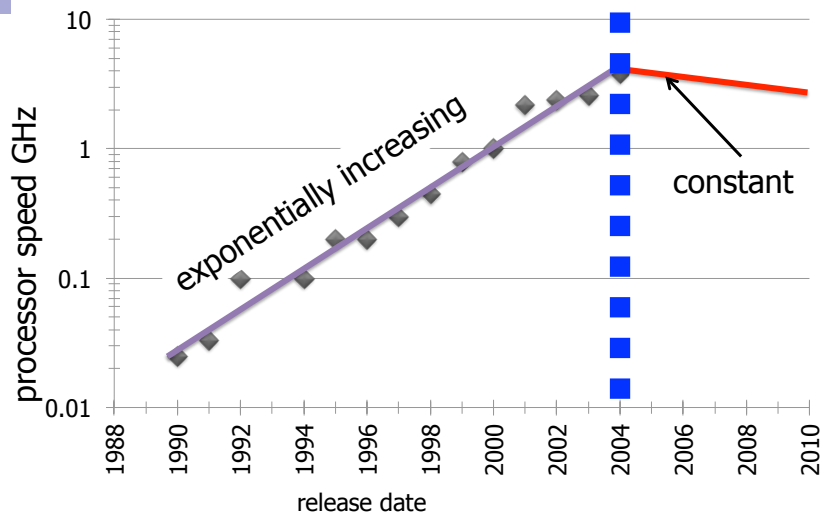
The Age of Big Data

By STEVE LOHR

Published: February 11, 2012

“... data a new class of economic
asset, like currency or gold.”

CPU's Stopped Getting Faster...



3

ML in the Context of Parallel Architectures



- But scalable ML in these systems is hard, especially in terms of:
 1. Programmability
 2. Data distribution
 3. Failures

©Carlos Guestrin 2013

4

Programmability Challenge 1: Designing Parallel programs

- SGD for LR:

- For each data point $\mathbf{x}^{(t)}$:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + \phi_i(\mathbf{x}^{(t)}) [y^{(t)} - P(Y = 1 | \phi(\mathbf{x}^{(t)}), \mathbf{w}^{(t)})] \right\}$$

©Carlos Guestrin 2013

5

Programmability Challenge 2: Race Conditions

- We are used to sequential programs:

- Read data, think, write data, read data, think, write data, read data, think, write data, read data, think, write data, read data, think, write data, read data, think, write data...

- But, in parallel, you can have non-deterministic effects:

- One machine reading data while other is writing

- Called a race-condition:

- Very annoying
 - One of the hardest problems to debug in practice:
 - because of non-determinism, bugs are hard to reproduce

©Carlos Guestrin 2013

6

Data Distribution Challenge

- Accessing data:
 - Main memory reference: 100ns (10^{-7} s)
 - Round trip time within data center: 500,000ns (5×10^{-4} s)
 - Disk seek: 10,000,000ns (10^{-2} s)
- Reading 1MB sequentially:
 - Local memory: 250,000ns (2.5×10^{-4} s)
 - Network: 10,000,000ns (10^{-2} s)
 - Disk: 30,000,000ns (3×10^{-2} s)
- Conclusion: Reading data from local memory is **much** faster → Must have data locality:
 - Good data partitioning strategy fundamental!
 - “Bring computation to data” (rather than moving data around)

©Carlos Guestrin 2013

7

Robustness to Failures Challenge

- From Google’s Jeff Dean, about their clusters of 1800 servers, in first year of operation:
 - 1,000 individual machine failures
 - thousands of hard drive failures
 - one power distribution unit will fail, bringing down 500 to 1,000 machines for about 6 hours
 - 20 racks will fail, each time causing 40 to 80 machines to vanish from the network
 - 5 racks will “go wonky,” with half their network packets missing in action
 - the cluster will have to be rewired once, affecting 5 percent of the machines at any given moment over a 2-day span
 - 50% chance cluster will overheat, taking down most of the servers in less than 5 minutes and taking 1 to 2 days to recover
- How do we design distributed algorithms and systems robust to failures?
 - It’s not enough to say: run, if there is a failure, do it again... because you may never finish

©Carlos Guestrin 2013

8

Move Towards Higher-Level Abstraction

- Distributed computing challenges are hard and annoying!
 1. Programmability
 2. Data distribution
 3. Failures
- High-level abstractions try to simplify distributed programming by hiding challenges:
 - Provide different levels of robustness to failures, optimizing data movement and communication, protect against race conditions...
 - Generally, you are still on your own WRT designing parallel algorithms
- Some common parallel abstractions:
 - Lower-level:
 - Pthreads: abstraction for distributed threads on single machine
 - MPI: abstraction for distributed communication in a cluster of computers
 - Higher-level:
 - Map-Reduce (Hadoop: open-source version): mostly data-parallel problems
 - GraphLab: for graph-structured distributed problems

©Carlos Guestrin 2013

9

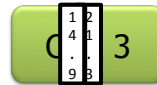
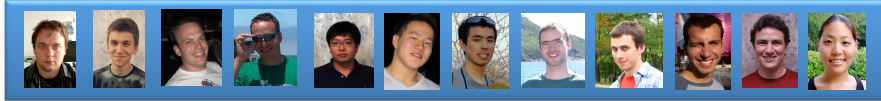
Simplest Type of Parallelism: Data Parallel Problems

- You have already learned a classifier
 - What's the test error?
- You have 10B labeled documents and 1000 machines
- Problems that can be broken into independent subproblems are called data-parallel (or embarrassingly parallel)
- Map-Reduce is a great tool for this...
 - Focus of today's lecture
 - but first a simple example

©Carlos Guestrin 2013

10

Data Parallelism (MapReduce)



*Solve a huge number of **independent** subproblems,
e.g., extract features in images*

Counting Words on a Single Processor



- (This is the “Hello World!” of Map-Reduce)
- Suppose you have 10B documents and 1 machine
- You want to count the number of appearances of each word on this corpus
 - Similar ideas useful, e.g., for building Naïve Bayes classifiers and computing TF-IDF
- Code:

Naïve Parallel Word Counting

- Simple data parallelism approach:
- Merging hash tables: annoying, potentially not parallel → no gain from parallelism???

©Carlos Guestrin 2013

13

Counting Words in Parallel & Merging Hash Tables in Parallel

- Generate pairs (word,count)
- Merge counts for each word in parallel
 - Thus parallel merging hash tables

©Carlos Guestrin 2013

14

Map-Reduce Abstraction

- Map:
 - Data-parallel over elements, e.g., documents
 - Generate (key,value) pairs
 - "value" can be any data type

- Reduce:
 - Aggregate values for each key
 - Must be commutative-associate operation
 - Data-parallel over keys
 - Generate (key,value) pairs

- Map-Reduce has long history in functional programming
 - But popularized by Google, and subsequently by open-source Hadoop implementation from Yahoo!

©Carlos Guestrin 2013

15

Map Code (Hadoop): Word Count

```
public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(LongWritable key, Text value, Context context) throws <stuff>
    {
        String line = value.toString();
        StringTokenizer tokenizer = new StringTokenizer(line);
        while (tokenizer.hasMoreTokens()) {
            word.set(tokenizer.nextToken());
            context.write(word, one);
        }
    }
}
```

©Carlos Guestrin 2013

16

Reduce Code (Hadoop): Word Count

```
public static class Reduce extends Reducer<Text, IntWritable,
Text, IntWritable> {

    public void reduce(Text key, Iterable<IntWritable> values,
        Context context)
        throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        context.write(key, new IntWritable(sum));
    }
}
```

©Carlos Guestrin 2013

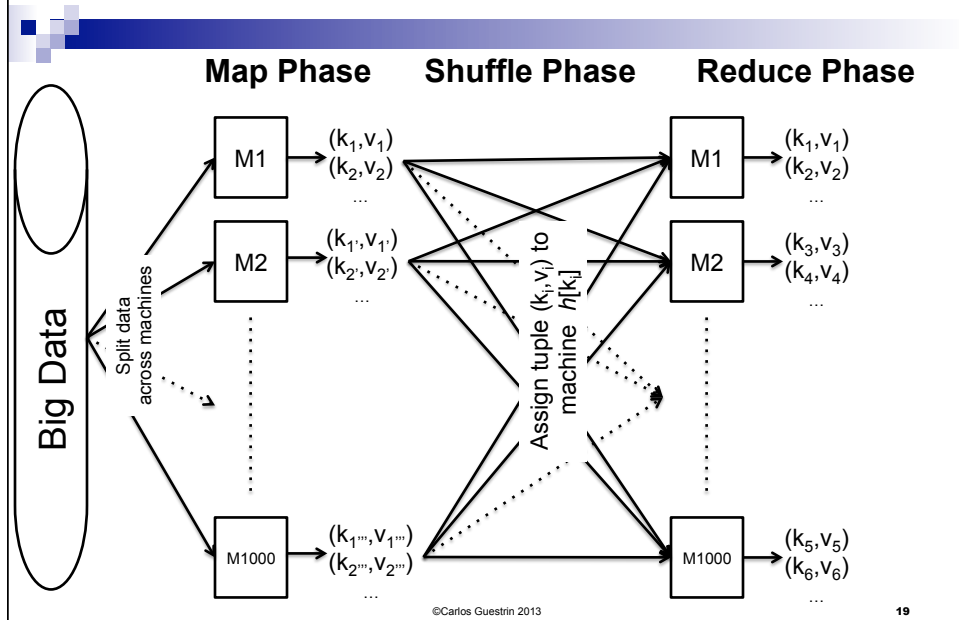
17

Map-Reduce Parallel Execution

©Carlos Guestrin 2013

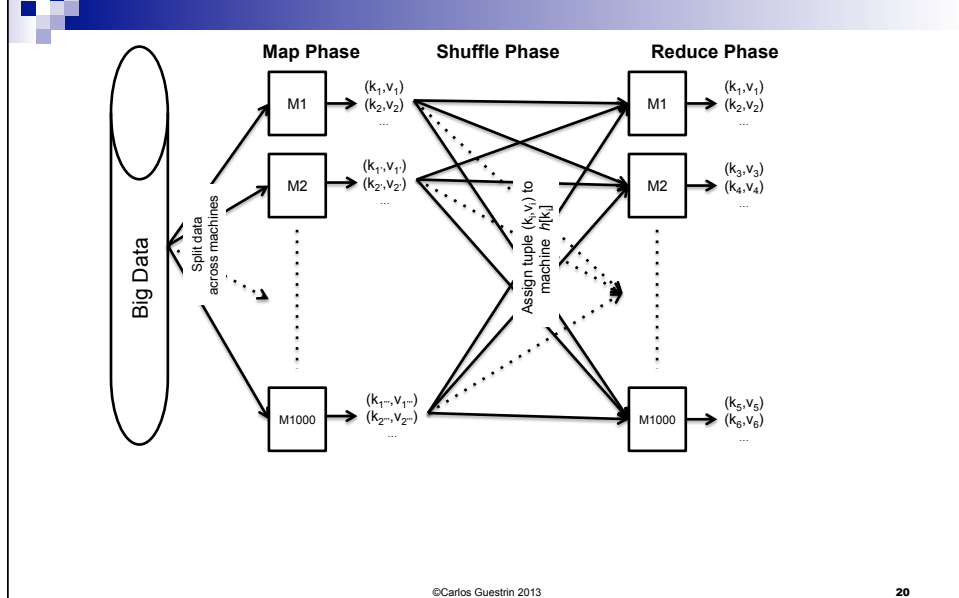
18

Map-Reduce – Execution Overview



19

Map-Reduce – Robustness to Failures 1: Protecting Data: **Save To Disk Constantly**



20

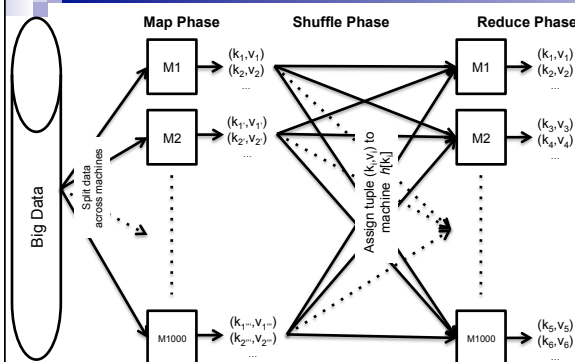
Distributed File Systems

- Saving to disk locally is not enough → If disk or machine fails, all data is lost
- Replicate data among multiple machines!
- Distributed File System (DFS)
 - Write a file anywhere → automatically replicated
 - Can read a file anywhere → read from closest copy
 - If failure, try next closest copy
- Common implementations:
 - Google File System (GFS)
 - Hadoop File System (HDFS)
- Important practical considerations:
 - Write large files
 - Many small files → becomes way too slow
 - Typically, files can't be "modified", just "replaced" → makes robustness much simpler

©Carlos Guestrin 2013

21

Map-Reduce – Robustness to Failures 2: Recovering From Failures: **Read from DFS**



- Communication in initial distribution & shuffle phase "automatic"
 - Done by DFS
- If failure, don't restart everything
 - Otherwise, never finish
- Only restart Map/Reduce jobs in dead machines

©Carlos Guestrin 2013

22

Improving Performance: Combiners

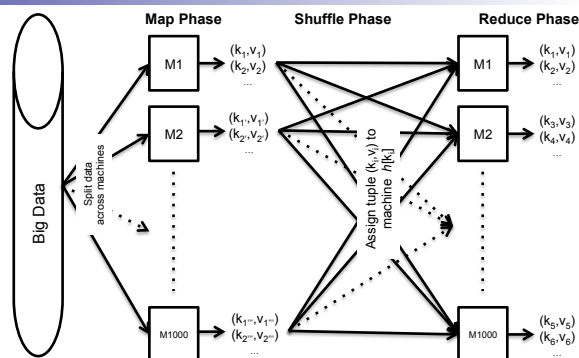
- Naïve implementation of M-R very wasteful in communication during shuffle:
- **Combiner:** Simple solution, perform reduce locally before communicating for global reduce
 - Works because reduce is commutative-associative

©Carlos Guestrin 2013

23

(A few of the) Limitations of Map-Reduce

- Too much synchrony
 - E.g., reducers don't start until all mappers are done
- “Too much” robustness
 - Writing to disk all the time
- Not all problems fit in Map-Reduce
 - E.g., you can't communicate between mappers
- Oblivious to structure in data
 - E.g., if data is a graph, can be much more efficient
 - For example, no need to shuffle nearly as much
- Nonetheless, extremely useful; industry standard for Big Data
 - Though many many companies are moving away from Map-Reduce (Hadoop)



©Carlos Guestrin 2013

24

What you need to know about Map-Reduce

- Distributed computing challenges are hard and annoying!
 1. Programmability
 2. Data distribution
 3. Failures
- High-level abstractions help a lot!
- Data-parallel problems & Map-Reduce
- Map:
 - Data-parallel transformation of data
 - Parallel over data points
- Reduce:
 - Data-parallel aggregation of data
 - Parallel over keys
- Combiner helps reduce communication
- Distributed execution of Map-Reduce:
 - Map, shuffle, reduce
 - Robustness to failure by writing to disk
 - Distributed File Systems

©Carlos Guestrin 2013

25

Case Study 2: Document Retrieval

Parallel K-Means on Map-Reduce

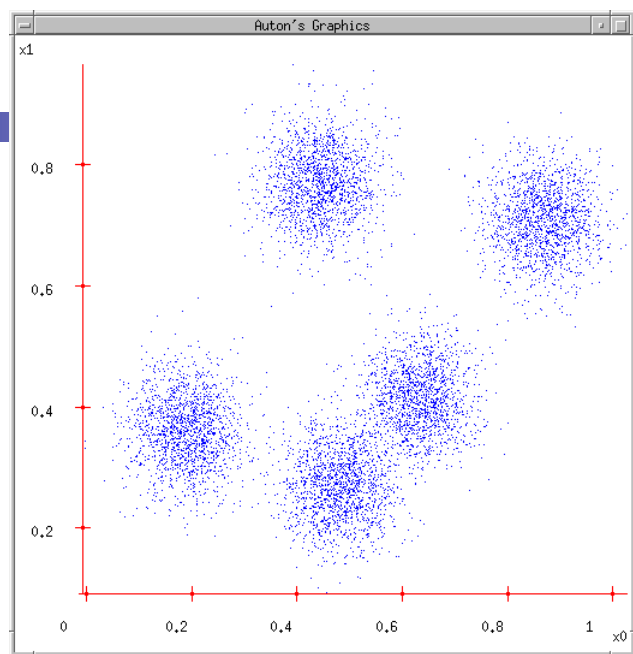
Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin
January 31st, 2013

©Carlos Guestrin 2013

26

Some Data

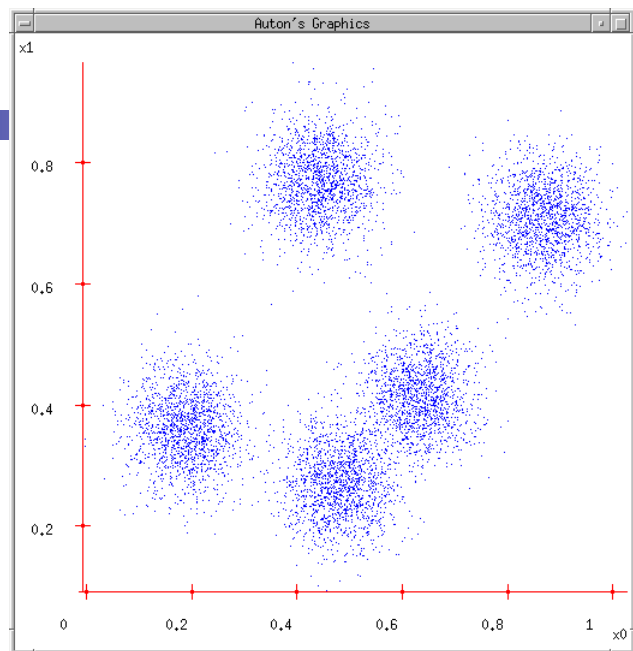


©2005-2009 Carlos Guestrin

27

K-means

1. Ask user how many clusters they'd like.
(e.g. $k=5$)

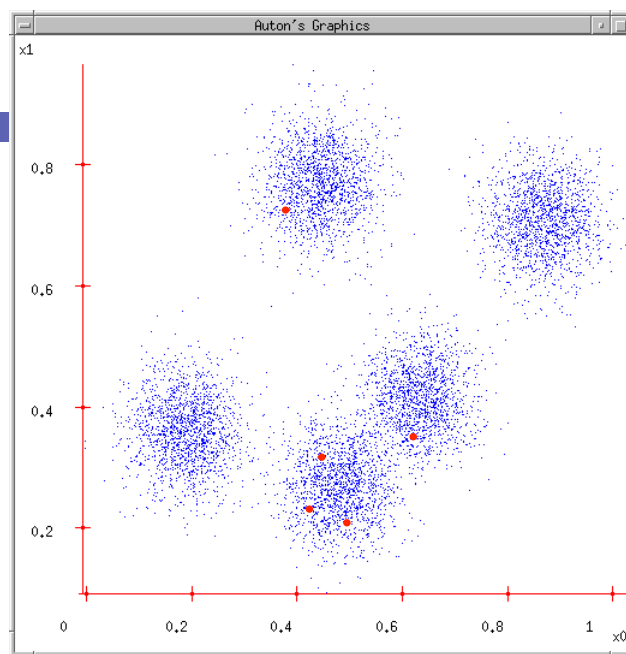


©2005-2009 Carlos Guestrin

28

K-means

1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations

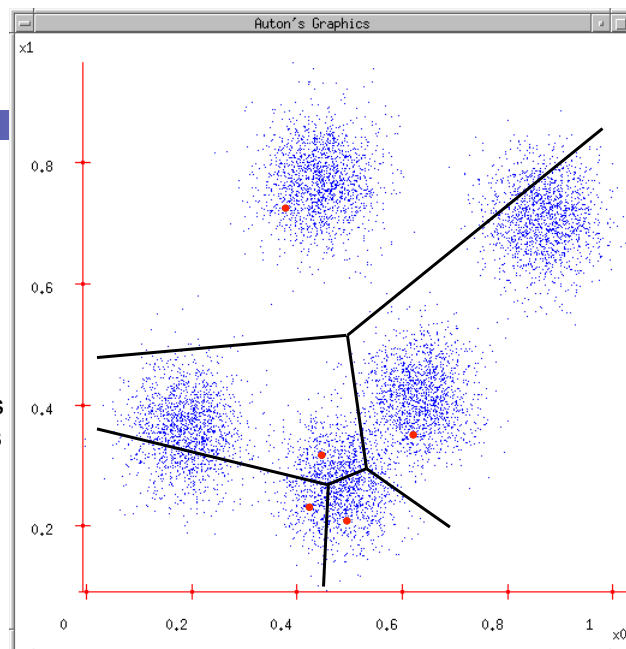


©2005-2009 Carlos Guestrin

29

K-means

1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)

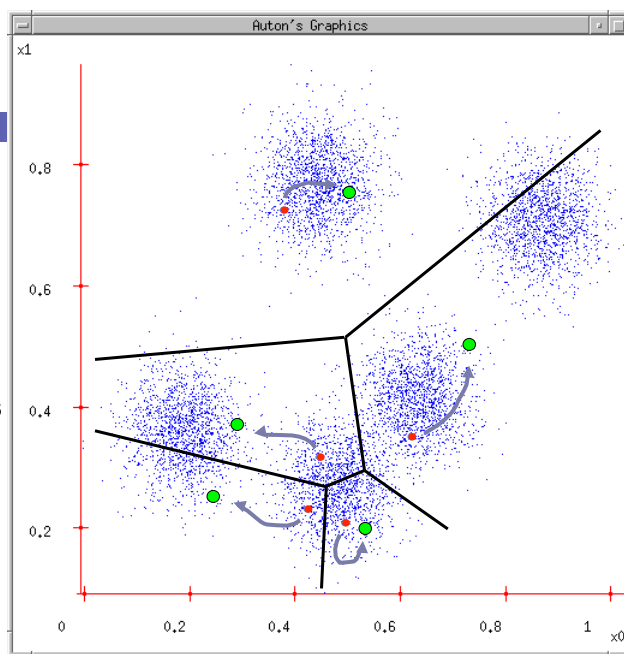


©2005-2009 Carlos Guestrin

30

K-means

1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns

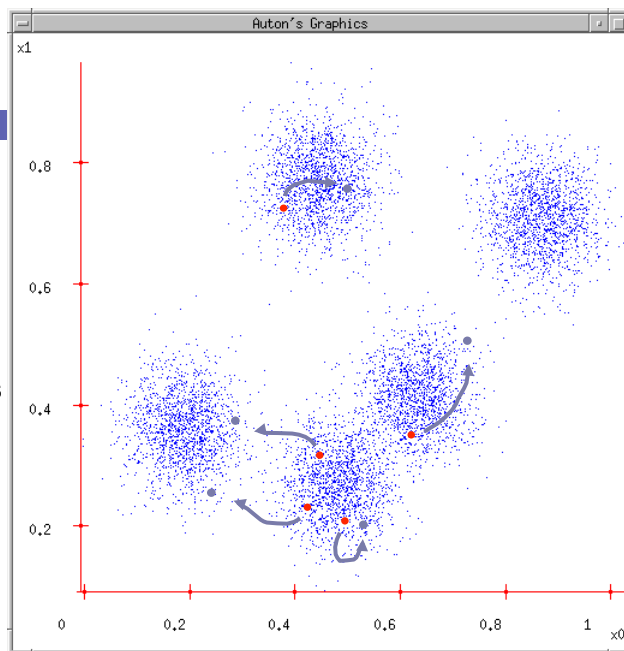


©2005-2009 Carlos Guestrin

31

K-means

1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



©2005-2009 Carlos Guestrin

32

K-means

- Randomly initialize k centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2$

- **Recenter:** μ_i becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} \|\mu - \mathbf{x}^j\|_2^2$

- Equivalent to $\mu_i \leftarrow$ average of its points!

©2005-2009 Carlos Guestrin

33

Special case: spherical Gaussians Mixtures and hard assignments

$$P(z = i | \mathbf{x}^j) \propto \frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x}^j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}^j - \mu_i) \right] P(z = i)$$

- If $P(Z=i|X)$ is spherical, with same σ for all classes:

$$P(z = i | \mathbf{x}^j) \propto \exp \left[-\frac{1}{2\sigma^2} \|\mathbf{x}^j - \mu_i\|^2 \right]$$

- If each \mathbf{x}^j belongs to one class z^j (hard assignment), marginal likelihood:

$$\prod_{j=1}^m \sum_{i=1}^k P(\mathbf{x}^j, z = i) \propto \prod_{j=1}^m \exp \left[-\frac{1}{2\sigma^2} \|\mathbf{x}^j - \mu_{z^j}\|^2 \right]$$

- Same as K-means!!!

©2005-2009 Carlos Guestrin

34

Map-Reducing One Iteration of K-Means

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $z^j \leftarrow \arg \min_i ||\mu_i - \mathbf{x}^j||_2^2$

- **Recenter:** μ_i becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} ||\mu - \mathbf{x}^j||_2^2$

- Equivalent to $\mu_i \leftarrow$ average of its points!

- **Map:**

- **Reduce:**

©2005-2009 Carlos Guestrin

35

Classification Step as Map

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $z^j \leftarrow \arg \min_i ||\mu_i - \mathbf{x}^j||_2^2$

- **Map:**

©Carlos Guestrin 2013

36

Recenter Step as Reduce

- **Recenter:** μ_i becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} \|\mu - \mathbf{x}^j\|_2^2$

- Equivalent to $\mu_i \leftarrow$ average of its points!

- **Reduce:**

©Carlos Guestrin 2013

37

Some Practical Considerations

- K-Means needs an iterative version of Map-Reduce
 - Not standard formulation
- Mapper needs to get data point and all centers
 - A lot of data!
 - Better implementation: mapper gets many data points

©Carlos Guestrin 2013

38

What you need to know about Parallel K-Means on Map-Reduce

- K-Means = EM for mixtures of spherical Gaussians with hard assignments
- Map: classification step; data parallel over data point
- Reduce: recompute means; data parallel over centers