**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

**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

**Sunday Review**

WORLD    U.S.    N.Y. / REGION    BUSINESS    TEC
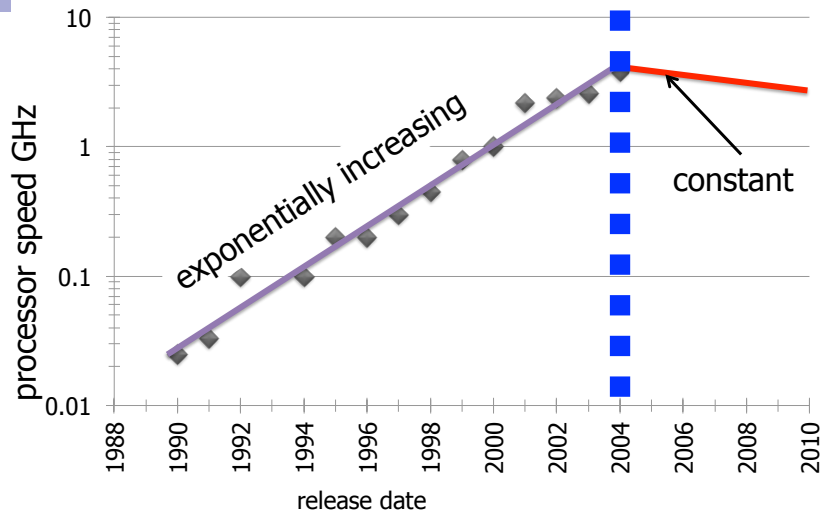
NEWS ANALYSIS
The Age of Big Data
By STEVE LOHR
Published: February 11, 2012

"… data a new class of economic asset, like currency or gold."

# CPUs Stopped Getting Faster…



# ML in the Context of Parallel Architectures



GPUs     Multicore     Clusters     Clouds     Supercomputers

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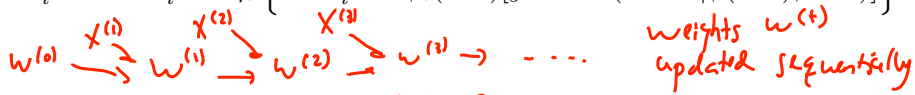  1. Programmability
  2. Data distribution
  3. Failures

4

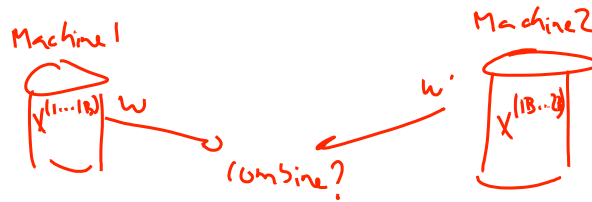# Programmability Challenge 1: Designing Parallel programs

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

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + \phi_i(\mathbf{x}^{(t)})[y^{(t)} - P(Y=1|\phi(\mathbf{x}^{(t)}), \mathbf{w}^{(t)})] \right\}$$
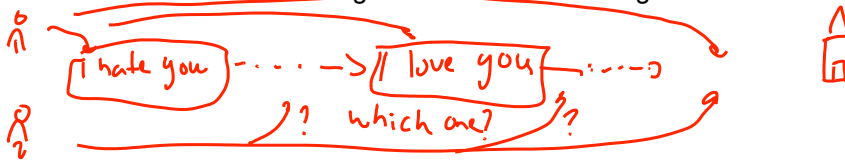
5

# 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 will other is writing
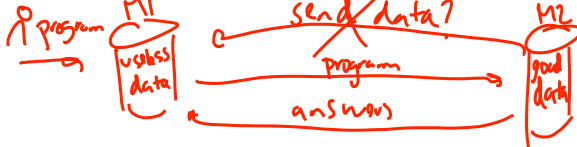


- Called a race-condition:
  - Very annoying
  - One of the hardest problems to debug in practice:
    - because of non-determinism, bugs are hard to reproduce

6

3

# Data Distribution Challenge

- Accessing data:
  - Main memory reference: 100ns ($10^{-7}$s) ← *local data*
  - Round trip time within data center: 500,000ns ($5 * 10^{-4}$s) ← *net access*
  - Disk seek: 10,000,000ns ($10^{-2}$s)
- Reading 1MB sequentially:
  - Local memory: 250,000ns ($2.5 * 10^{-4}$s) *} 2 orders of magnitude*
  - 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)

*[handwritten diagram: program → M1 "useless data", send data?, program, answers, M2 "good data"]*

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
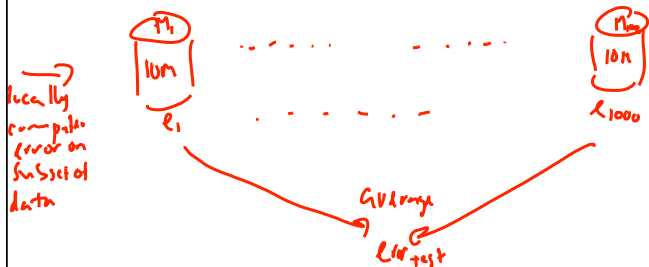
8

4

## 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  *(handwritten: this quarter)*
    - GraphLab: for graph-structured distributed problems

©Carlos Guestrin 2013    9

---

## Simplest Type of Parallelism: Data Parallel Problems

- You have already learned a classifier $w^*$
  - □ What's the test error?

$$\ell_{rr} = \frac{1}{N_{test}} \sum_i \frac{1}{2} | y^{(i)} - sign(w^* \cdot x^{(i)}) |$$

- You have 10B labeled documents and 1000 machines

*(handwritten diagram with: locally compute error on subset of data; cylinders labeled m_1, 10m, ℓ_1; m, 10n, ℓ_1000; average, ℓ_0 test)*

- 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)



CPU 1    CPU 2    CPU 3    CPU 4

$x^{(1)}$   $x^{(2)}$      $x^{(12)}$

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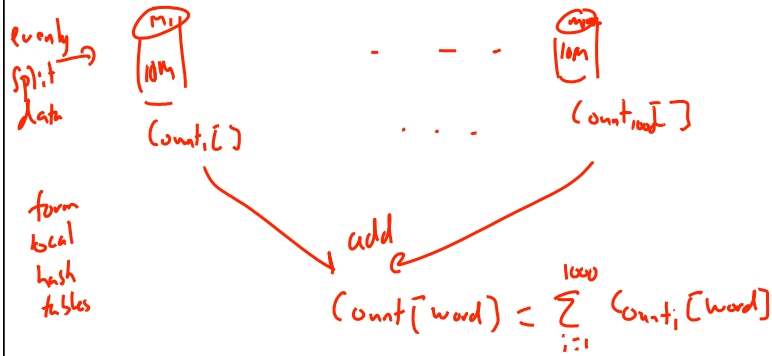e.g., extract features in images**

---

# 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:

```
Count[] ← int a hash table
for d in documents
    for word in d
        Count[word] += 1
```

©Carlos Guestrin 2013

12

6

# Naïve Parallel Word Counting

- Simple data parallelism approach:

evenly split data

form local hash tables

$M1$ $10M$ ... $10M$

$Count_1[\ ]$ ... $Count_{100}[\ ]$
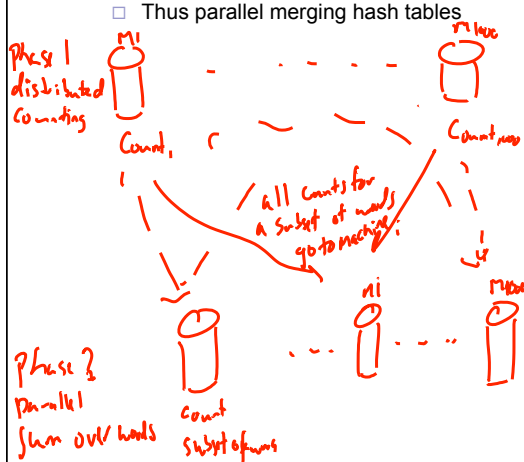
add

$$\left(Count[word] = \sum_{i=1}^{1000} Count_i[word]\right)$$

- Merging hash tables: annoying, potentially not parallel ➔ no gain from parallelism???

13

---

# Counting Words in Parallel & Merging Hash Tables in Parallel

- Generate pairs (word,count) ('UW', 17)
- Merge counts for each word in parallel
  - Thus parallel merging hash tables

Phase 1 distributed counting

$M1$ ... $M100$

$Count_1$ ... $Count_{100}$

all counts for a subset of words go to machine i

$M1$ $M100$

Phase 2 Parallel Sum over words

count subset of words

Which words go to machine:

$h: X \to [1, \ldots, \# \text{ machines}]$

send counts of word 'UW' to machine $h('UW')$

14

7

# Map-Reduce Abstraction

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

  *('UW', 17)*

  *in this example:* *('Mary', 1)*
  *('UW', 1)*
  *('Mary', 1)*

  *word count*
  *map (document)*
  *for word in doc*
  *emit (word, 1)*

- Reduce: *Take all values associated with a key and aggregate*
  - □ Aggregate values for each key
  - □ Must be commutative-associate operation
  - □ Data-parallel over keys
  - □ Generate (key,value) pairs

  *map reduce ('UW', [1, 17, 0, 0, 12])*
  *emit ('UW', 30)*

  *Reduce (word, count: list(int))*
  *c = 0*
  *for i in count*
  *c += count[i]*
  *emit (word, c)*

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

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();   ;doc

    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);
        }
    }
}
```

16

8

# 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));
    }
}
```
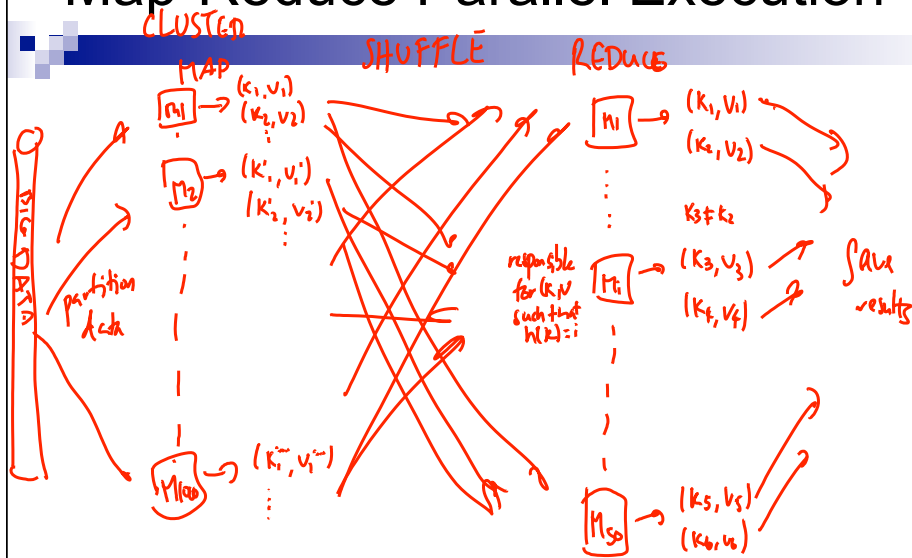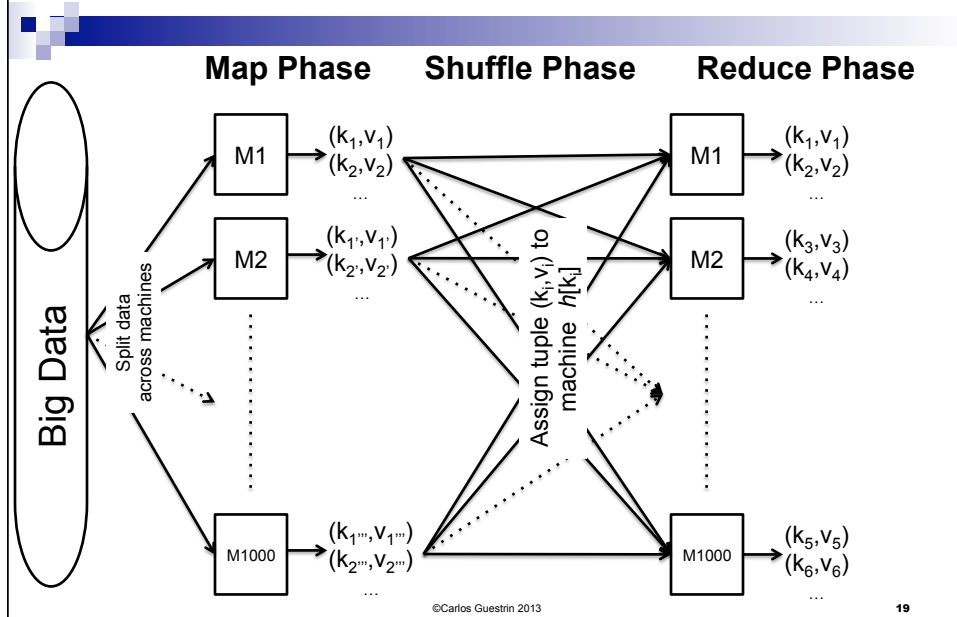
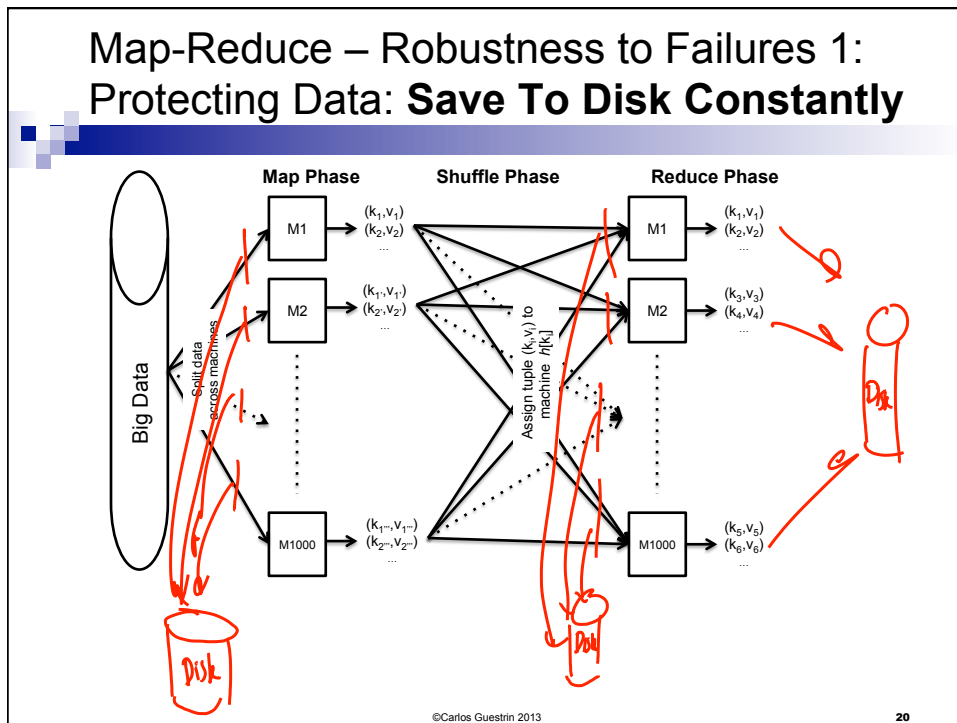*(handwritten annotations: word ✓ ; ✓ list of value ; Nr. each/count ; ← emit (word, total))*

17

# Map-Reduce Parallel Execution

18

9

Map-Reduce – Execution Overview


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

# Distributed File Systems

- Saving to disk locally is not enough → If disk or machine fails, all data is lost
- Replicate data among multiple machines!  *Usually 3*
- 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:
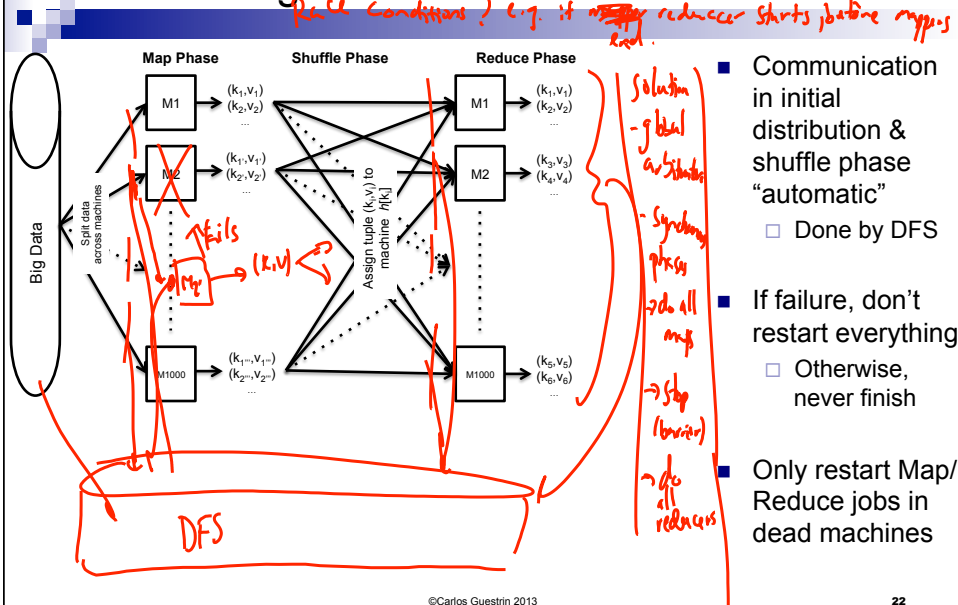  - Google File System (GFS)
  - Hadoop File System (HDFS)
- Important practical considerations:
  - Write large files
    - Many small files → becomes way too slow
  - Typically, files can't be "modified", just "replaced" → makes robustness much simpler

©Carlos Guestrin 2013
21

---

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

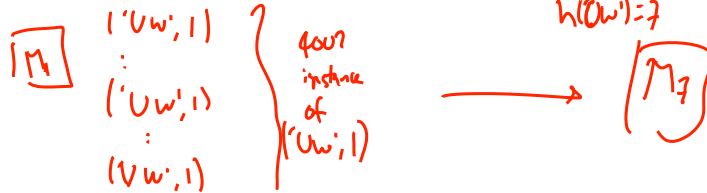**Map Phase**    **Shuffle Phase**    **Reduce Phase**

- 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
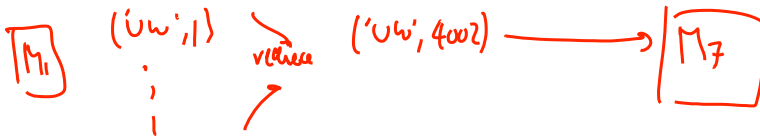
©Carlos Guestrin 2013
22

11

# Improving Performance: Combiners

■ Naïve implementation of M-R very wasteful in communication during shuffle:

*(handwritten annotations)*

$h('Uw')=7$

four instead of $('Uw', 1)$

■ **Combiner**: Simple solution, perform reduce locally before communicating for global reduce
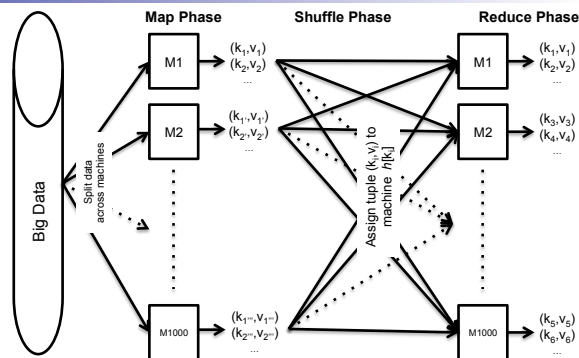
  □ Works because reduce is commutative-associative

*(handwritten: $('Uw', 1) \xrightarrow{reduce} ('Uw', 4002) \rightarrow M_7$)*

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)



**Map Phase**    **Shuffle Phase**    **Reduce Phase**

24

12

# 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

---

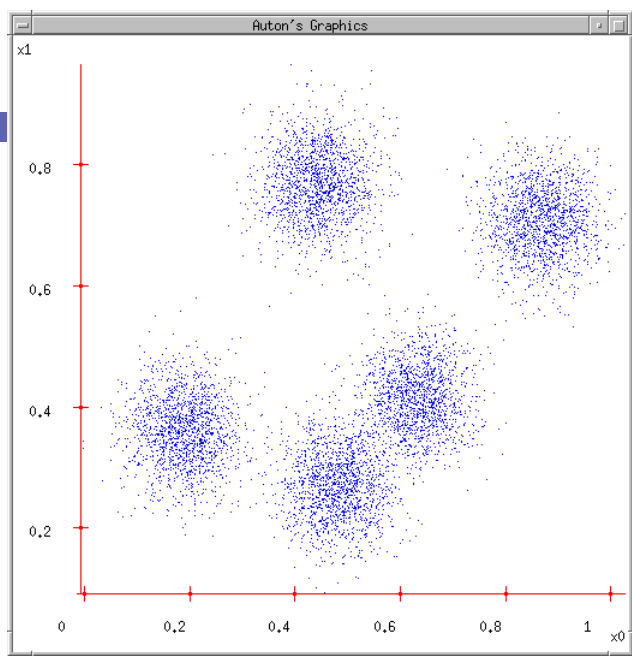# Case Study 2: Document Retrieval

# Parallel K-Means on Map-Reduce

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

Carlos Guestrin
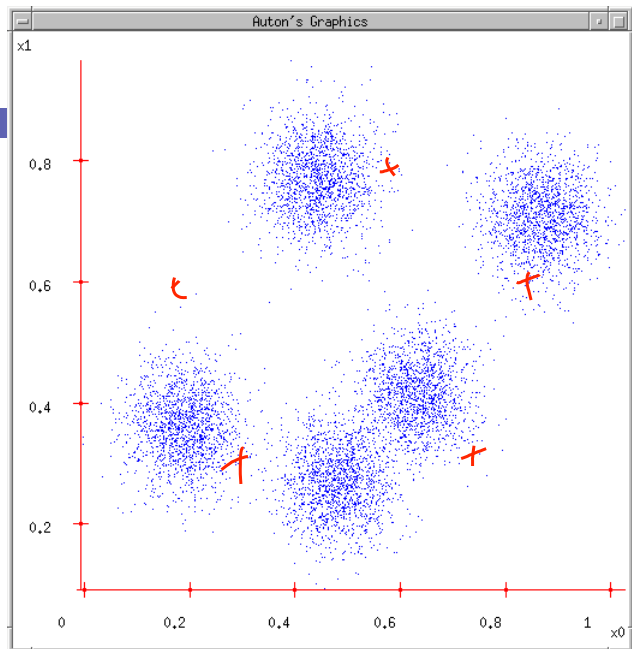
January 31st, 2013

©Carlos Guestrin 2013

26

# Some Data

27

# K-means

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

28

# K-means

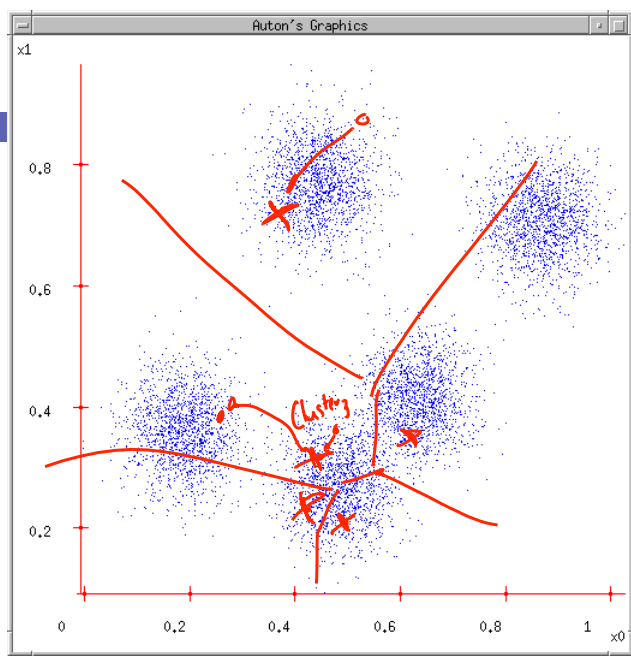1. Ask user how many clusters they'd like. *(e.g. k=5)*

2. Randomly guess k cluster Center locations

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)

30

15
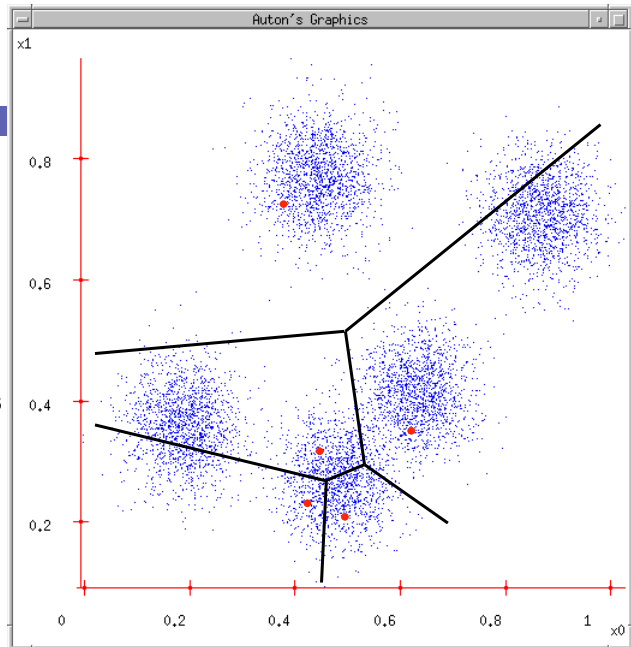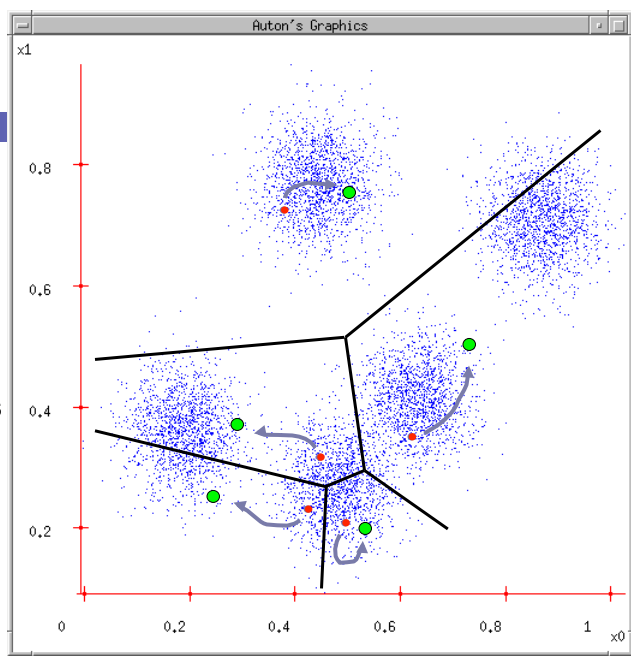
# K-means

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

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.
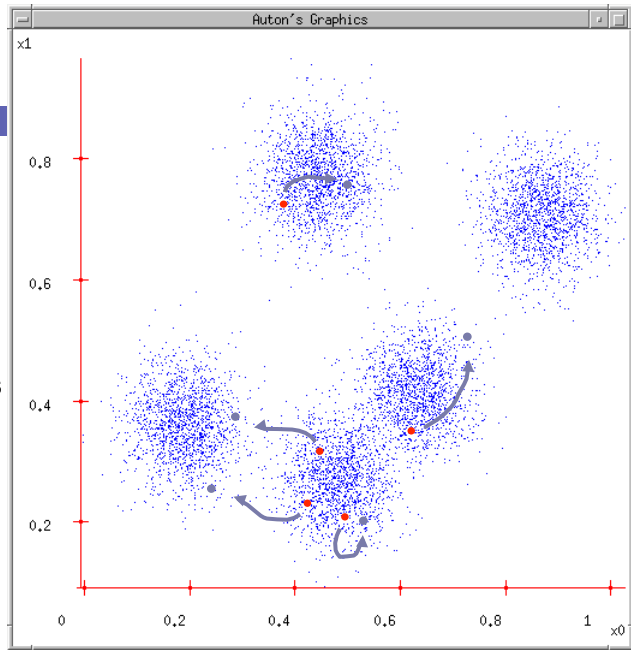
4. Each Center finds the centroid of the points it owns…

5. …and jumps there

6. …Repeat until terminated!

32

16

# K-means

- Randomly initialize $k$ centers
  - $\mu^{(0)} = \mu_1^{(0)}, \ldots, \mu_k^{(0)}$

- **Classify**: Assign each point $j \in \{1, \ldots m\}$ to nearest center:
  *(assign points to nearest cluster center)*
  - $z^j \leftarrow \arg\min_i ||\mu_i - \mathbf{x}^j||_2^2$

- **Recenter**: $\mu_i$ 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!

$$\mu_i = \frac{\sum_{j:z^j=i} x^j}{\sum_{j:z^j=i} 1}$$

33

---

# Special case: spherical Gaussians Mixtures and hard assignments $\Sigma = \sigma^2 I$

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

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

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

*K-means just like EM, but take maximal assignment*
*"E-step": $z^j = \arg\max_i P(z=i|x^j)$*

- If each $x^j$ belongs to one class $z^j$ (hard assignment), marginal likelihood:

$$\prod_{j=1}^m \sum_{i=1}^k P(\mathbf{x}^j, z=i) \propto \prod_{j=1}^m \exp\left[ -\frac{1}{2\sigma^2}\left\|\mathbf{x}^j - \mu_{z^j}\right\|^2 \right] = \frac{-1}{2\sigma^2} \sum_{j=1}^m \|x^j - \mu_{z^j}\|_2^2$$

- Same as K-means!!!    *if only solve for $\mu$*

34

17

# Map-Reducing One Iteration of K-Means

- **Classify**: Assign each point $j \in \{1,...m\}$ to nearest center:
  - $z^j \leftarrow \arg\min_i ||\mu_i - \mathbf{x}^j||_2^2$

- **Recenter**: $\mu_i$ 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**: data-parallel: classify phase
  for each data point: given $(\mu, \mathbf{x})$ → emit $(z^j, \mathbf{x}^j)$

- **Reduce**: Recenter phase: average over all points in class $i$

35

---

# Classification Step as Map

- **Classify**: Assign each point $j \in \{1,...m\}$ to nearest center:
  - $z^j \leftarrow \arg\min_i ||\mu_i - \mathbf{x}^j||_2^2$

- **Map**: map $([\mu_1,...,\mu_K], \mathbf{x}^j)$

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

  emit $(z^j, \mathbf{x}^j)$

  $z^j = 2$
  e.g. emit$(2, (17, op, 1, 7, 2))$

36

# Recenter Step as Reduce

■ **Recenter**: $\mu_i$ 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**: Reduce( i , list_x : $[x^1, x^2, \ldots]$ )

*which were assigned to class i*

  Count = 0
  Sum = 0
  for x in list_x
      sum += x
      count += 1
  Emit ( i , sum/count )

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

38

19

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

- K-Means = EM for mixtures of spherical Gaussians with hard assignments $in$ $E-Stp$

- Map: classification step; data parallel over data point

- Reduce: recompute means; data parallel over centers