**Case Study 2: Document Retrieval**

# Clustering Documents

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

January 23rd, 2014

**1**

---

# Document Retrieval

- **Goal:** Retrieve documents of interest
- **Challenges:**
  - ☐ Tons of articles out there
  - ☐ How should we measure similarity?

ARTICLES

**2**

# Task 1: Find Similar Documents

■ **So far…**

☐ **Input:** Query article $x$

☐ **Output:** Set of k similar articles

*k-NN approach*

$\{x^{NN_1}, \ldots, x^{NN_k}\}$

# Task 2: Cluster Documents

■ **Now:**

☐ Cluster documents based on topic

"sports"

"World news"

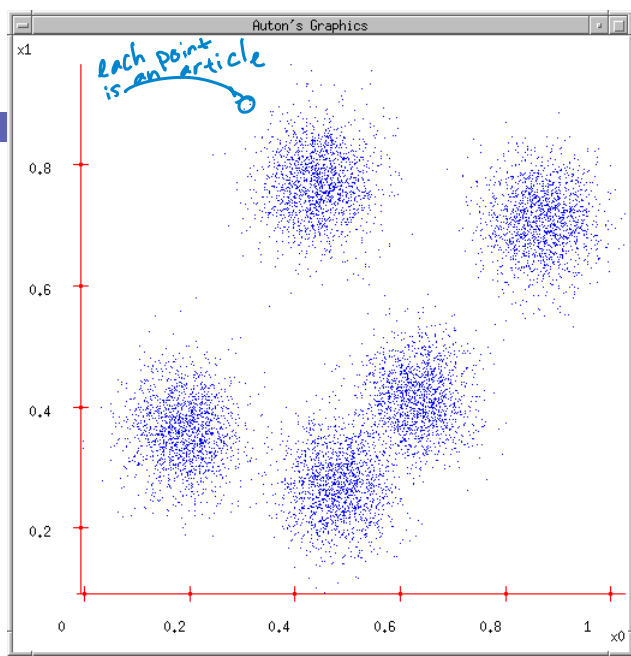# Some Data

e.g. doc

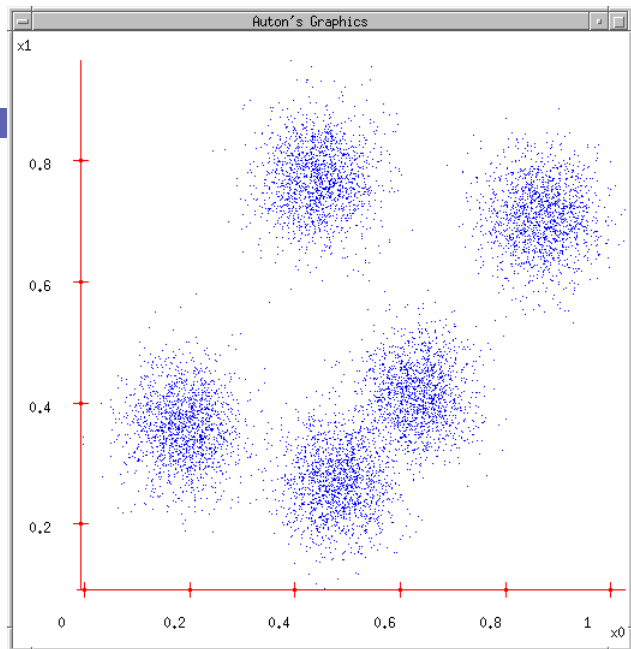X bag of words
or
tfidf

How to discover
clusters?

(unsupervised)

each point
is an article

x1

0.8

0.6

0.4

0.2

0          0.2        0.4        0.6        0.8        1    x0

Auton's Graphics

©Emily Fox 2014                    5

# K-means

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

x1

0.8

0.6

0.4

0.2

0          0.2        0.4        0.6        0.8        1    x0

Auton's Graphics
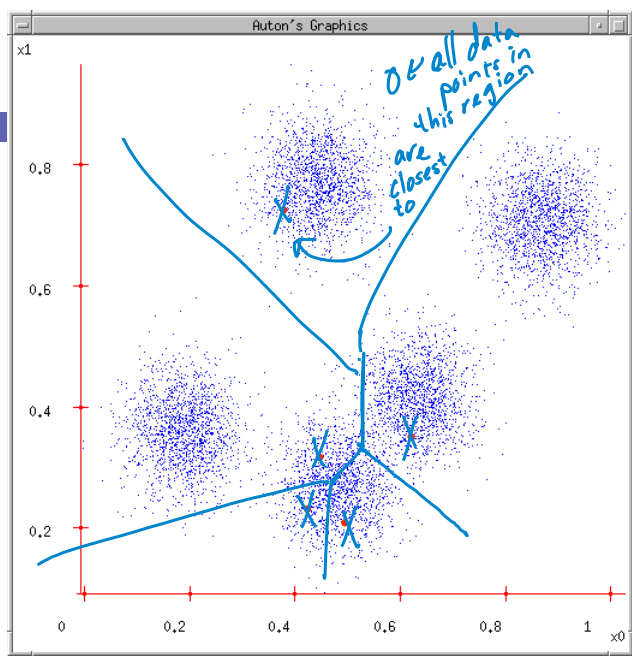
©Emily Fox 2014                    6

3

# K-means

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

2. Randomly guess k cluster Center locations

---

# 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)
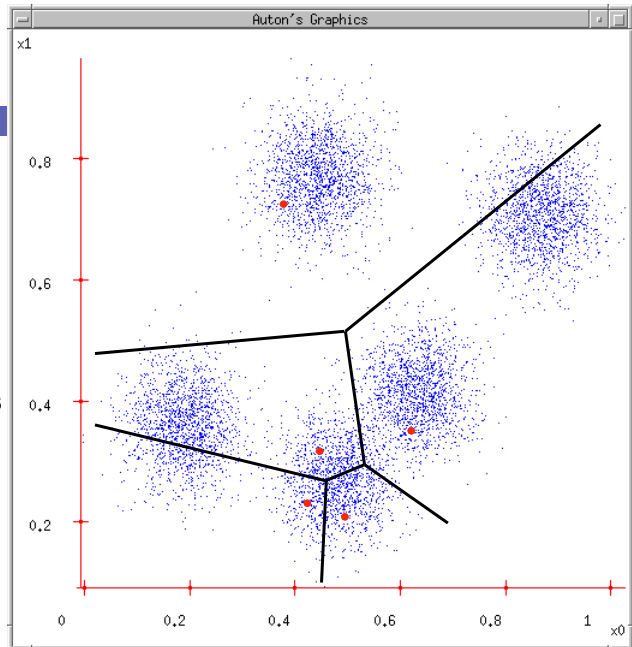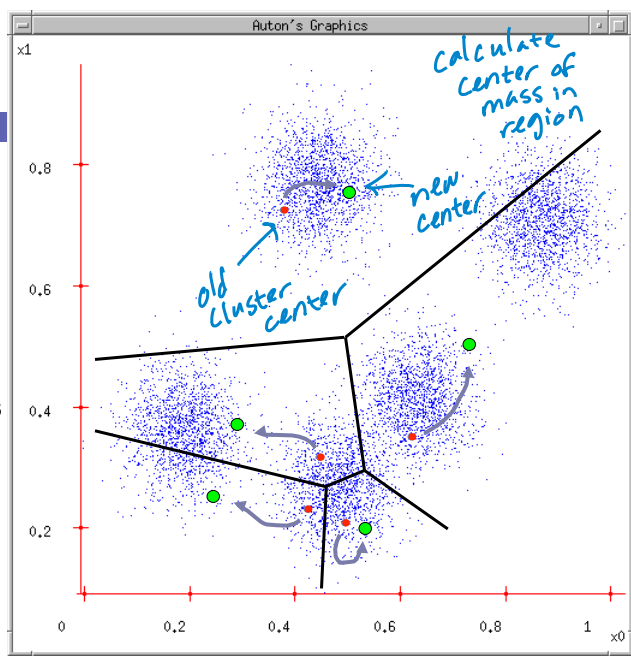
4

# 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



*Auton's Graphics*

calculate center of mass in region

new center

old cluster center

9

---

# 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.

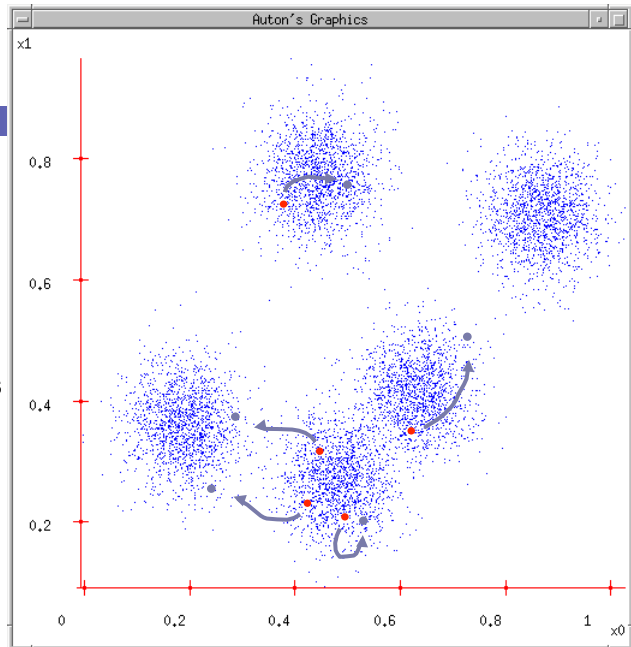4. Each Center finds the centroid of the points it owns…

5. …and jumps there

6. …Repeat until terminated!



*Auton's Graphics*

10

# K-means

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

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

  *z^j is indicator of nearest center ("cluster indicator")*

- **Recenter**: $\mu_i$ becomes centroid of its point:
  - □ $\mu_i^{(t+1)} \leftarrow \arg\min_\mu \sum_{j:z^j=i} ||\mu - \mathbf{x}^j||_2^2$

  $M_i = \sum_{j:z^j=i} x^j$

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

  *# of pts in cluster* $|\{x^j : z^j = i\}|$

©Emily Fox 2014

11

---

## Case Study 2: Document Retrieval

# Parallel Programming
# Map-Reduce

Machine Learning/Statistics for Big Data
CSE547/STAT548, University of Washington
Emily Fox
January 23rd, 2014

©Emily Fox 2014

12

# 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**

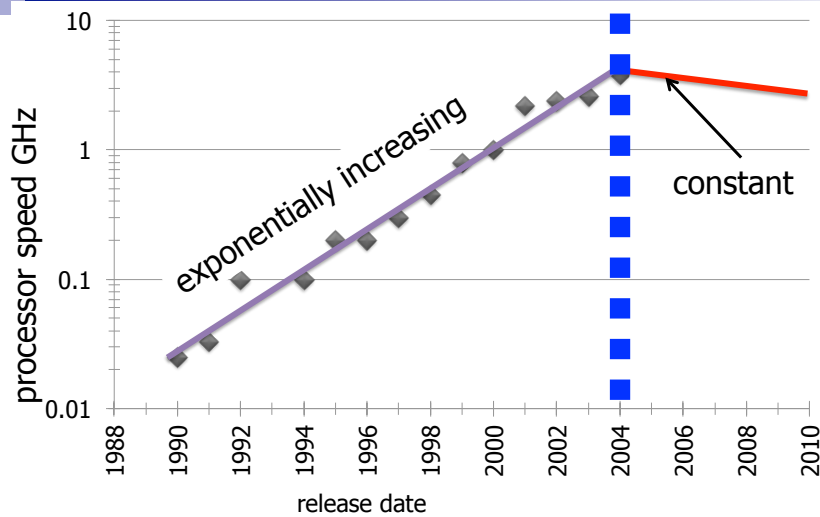WORLD    U.S.    N.Y. / REGION    BUSINESS    TEC

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…



- processor speed GHz (y-axis): 0.01, 0.1, 1, 10
- release date (x-axis): 1988, 1990, 1992, 1994, 1996, 1998, 2000, 2002, 2004, 2006, 2008, 2010
- *exponentially increasing*
- constant

14

# ML in the Context of Parallel Architectures

*use more processors*

GPUs     Multicore     Clusters     Clouds     Supercomputers

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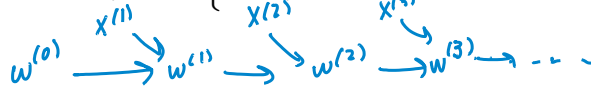  1. Programmability
  2. Data distribution
  3. Failures

*we'll go through these ideas...*

15

---

# Programmability Challenge 1: Designing Parallel programs

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

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

$x^{(1)}$   $x^{(2)}$   $x^{(3)}$

$w^{(0)} \longrightarrow w^{(1)} \longrightarrow w^{(2)} \longrightarrow w^{(3)} \longrightarrow \cdots$

*weights $w^{(t)}$ updated sequentially*

*How do we parallelize?*

*e.g. Machine 1*     *Machine 2*

$w$   $w'$

*Subset of data*   $(x^{(1)}, \ldots, x^{(B)})$   *combine?*   $(x^{(B+1)}, \ldots, x^{(2B)})$
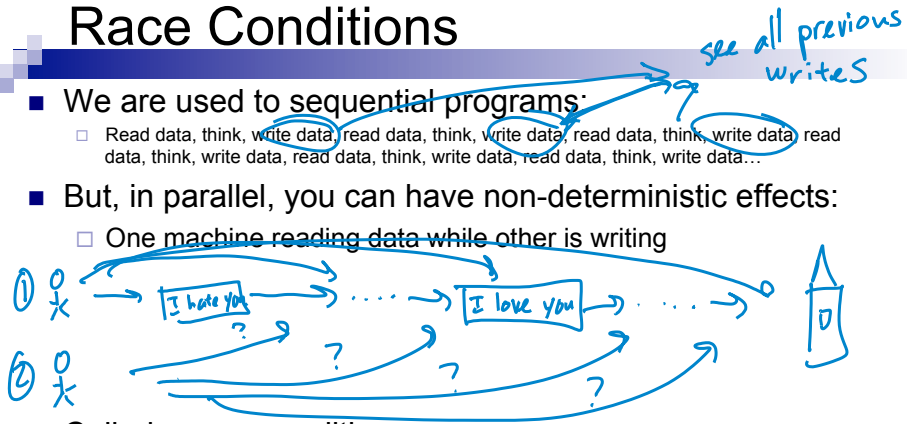
16

8

# 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:
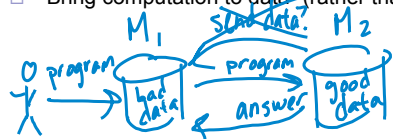    - because of non-determinism, bugs are hard to reproduce

17

# 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)

18

9

# 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

19

# 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
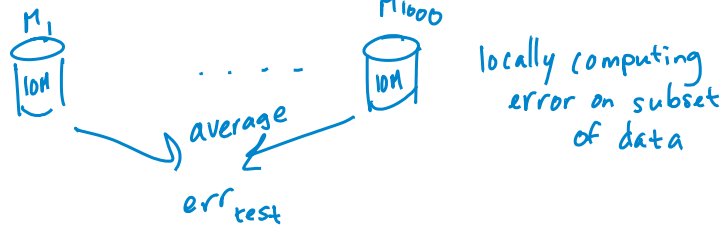    - GraphLab: for graph-structured distributed problems

  *this quarter*

20

10

# Simplest Type of Parallelism:
# Data Parallel Problems

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

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

- You have 10B labeled documents and 1000 machines

$M_1$     $M_{1000}$

10M     · · · · ·     10M

average

locally computing error on subset of data

$err_{test}$

- 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

21

# Data Parallelism (MapReduce)

CPU 1    1 7 . 5

CPU 2    6 7 . 5

CPU 3    1 4 . 9   2 1 . 3

CPU 4    3 4 . 3   2 5 . 8

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

22

# 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[ ]  ←——— init a hash table
for d in documents
    for word in d
        count[word] += 1
```

23

# Naïve Parallel Word Counting

- Simple data parallelism approach:

evenly split data

$n_1$   10M   . . . .   $n_{1000}$   10M

$Count_1[ ]$          $Count_{1000}[ ]$

add

$$Count[word] = \sum_{i=1}^{1000} Count_i[word]$$

How do we do this for all words in the vocab?

- Merging hash tables: annoying, potentially not parallel → no gain from parallelism??? have to merge hash tables sequentially

24

12

## 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*
*dist.*
*counting*

$M_1$      $M_{1000}$
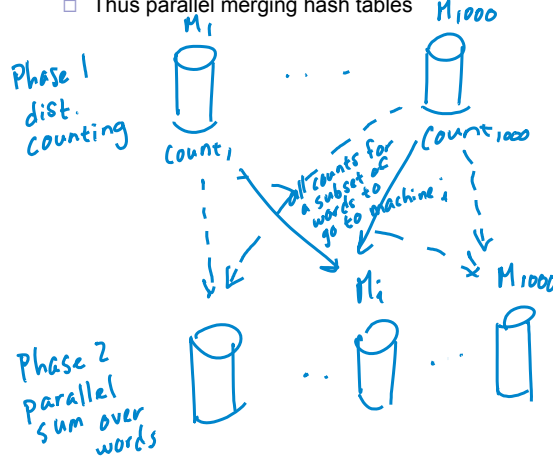
$Count_1$      $Count_{1000}$

*all counts for a subset of words to go to machine i*

$M_i$      $M_{1000}$

*Phase 2*
*parallel*
*sum over*
*words*

*which words*
*go to machine i:*

$h: \chi \rightarrow [1, .., \#\ machines]$

*send counts of*
*word 'UW'*
*to machine*

$h['UW']$

25

---

# 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)

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

reduce('UW', [1, 17, 0, 0, 12])

emit('UW', 30)

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

*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!

26

13

# Map Code (Hadoop): Word Count

*(handwritten: name)*  *(handwritten: class)*

```
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);
        }
    }
}
```

*(handwritten: for each word in String)*
*(handwritten: = word object)*
*(handwritten: emit (word, 1))*

27

# 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: word: 'UW')*
*(handwritten: to 0 count)*
*(handwritten: [1, 0, 0, 12])*
*(handwritten: in count)*
*(handwritten: emit (word, 30))*

28

14

# Map-Reduce Parallel Execution



©Emily Fox 2014    29

# Map-Reduce – Execution Overview



©Emily Fox 2014    30

15

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

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

Big Data

Split data across machines

M1 → $(k_1,v_1)$ $(k_2,v_2)$ …

M2 → $(k_1,v_1)$ $(k_2,v_2)$ …

M1000 → $(k_{1''},v_{1''})$ $(k_{2''},v_{2''})$ …

Assign tuple $(k_i,v_i)$ to machine $h[k_i]$

M1 → $(k_1,v_1)$ $(k_2,v_2)$ …

M2 → $(k_3,v_3)$ $(k_4,v_4)$ …

M1000 → $(k_5,v_5)$ $(k_6,v_6)$ …

disk

disk

disk

31

# 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 from anywhere ➔ automatically replicated
  - Can read a file from anywhere ➔ read from closest copy
    - If failure, try next closest copy

  *write('foo.txt')*
  *→ $h_i$ ['foo.txt]*
  *i=1,2,3*

  *get('foo.txt'): $h_i$('foo.txt') → where it is*
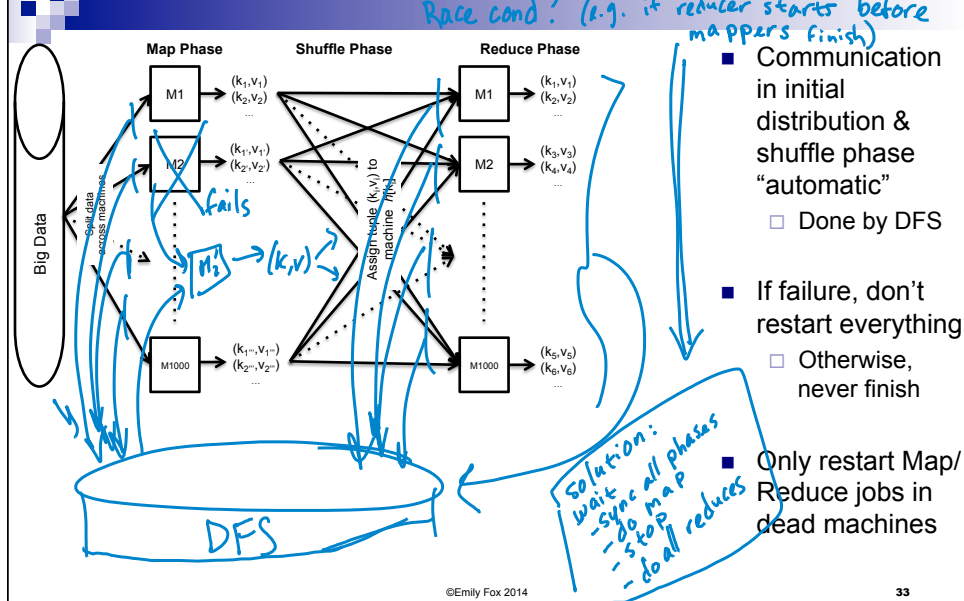  *failed → next one if failure*

- Common implementations:
  - Google File System (GFS)
  - Hadoop File System (HDFS)
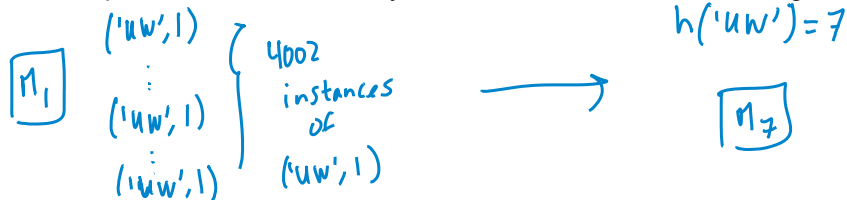- Important practical considerations:
  - Write large files
    - Many small files ➔ becomes way too slow
  - Typically, files can't be "modified", just "replaced" ➔ makes robustness much simpler

32

16

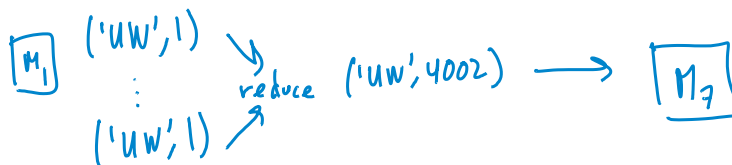## Map-Reduce – Robustness to Failures 2: Recovering From Failures: **Read from DFS**



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

*Race cond? (e.g. if reducer starts before mappers finish)*

Big Data

M1 → $(k_1,v_1)$ $(k_2,v_2)$ ...

M2 → $(k_1',v_1')$ $(k_2',v_2')$ ... *fails*

$M_2' \rightarrow (k,v)$

M1000 → $(k_1''',v_1''')$ $(k_2''',v_2''')$ ...

Assign tuple $(k,v)$ to machine $h(k)$

M1 → $(k_1,v_1)$ $(k_2,v_2)$ ...

M2 → $(k_3,v_3)$ $(k_4,v_4)$ ...

M1000 → $(k_5,v_5)$ $(k_6,v_6)$ ...

DFS

*Solution: wait all phases - sync all map - stop - do all reduces*

- ■ 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

©Emily Fox 2014    33

---

# Improving Performance: Combiners

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



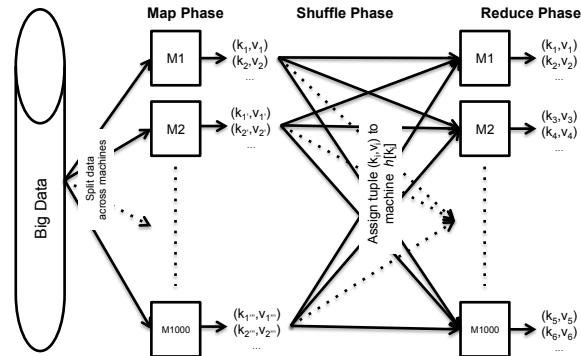$M_1$    ('uw',1) ... ('uw',1) ... ('uw',1)    4002 instances of ('uw',1) →    $h('uw')=7$    $M_7$

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

$M_1$    ('uw',1) ... ('uw',1) → reduce ('uw',4002) → $M_7$

©Emily Fox 2014    34

17

# (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**

Big Data — Split data across machines

M1 → $(k_1,v_1)$ $(k_2,v_2)$ …
M2 → $(k_1',v_1')$ $(k_2',v_2')$ …
M1000 → $(k_1''',v_1''')$ $(k_2''',v_2''')$ …

Assign tuple $(k_i, v_i)$ to machine $h[k_i]$

M1 → $(k_1,v_1)$ $(k_2,v_2)$ …
M2 → $(k_3,v_3)$ $(k_4,v_4)$ …
M1000 → $(k_5,v_5)$ $(k_6,v_6)$ …

35

---

# 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

36

18

**Case Study 2: Document Retrieval**

# Parallel K-Means on Map-Reduce

Machine Learning/Statistics for Big Data
CSE547/STAT548, University of Washington
Emily Fox
January 23$^{rd}$, 2014

37

---

# Map-Reducing One Iteration of K-Means

- **Classify**: Assign each point j∈{1,…*N*} 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 $(\{\mu_i\}, x^j) \rightarrow$ emit $(z^j, x^j)$
  key value

- **Reduce**: Recenter phase
  average over all points in class $i$   $z^j = i$

38

19

# Classification Step as Map

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

- **Map**:

$$map\left([\mu_1,\ldots,\mu_k], x^j\right)$$
$$z^j \leftarrow argmin \, ||\mu_i - x^j||_2^2$$
$$emit\left(z^j, x^j\right)$$

$$z^{j=2}$$
$$e.g. \quad emit\left(2, [17,0,0,1,7,2]\right)$$
$$x^j \text{ vector}$$

# 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!

  data assigned to class i

- **Reduce**:

$$Reduce\left(i, list\_x : [x^1, x^2, \ldots]\right)$$
$$count = 0$$
$$sum = 0$$
$$For \, x \, in \, list\_x$$
$$sum += x$$
$$count += 1$$
$$emit\left(i, \, sum/count\right)$$

avg. of pts in cluster i

# 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

*Each call to map processes multiple $x^j$'s ... faster*

41

---

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

- Map: classification step; data parallel over data point

- Reduce: recompute means; data parallel over centers

42