

2

L24: MapReduce and DFS

Distributed Systems for the masses

Jeff M. Phillips

April 15, 2019

Don't put too much, or small font!

- ▶ Succinct title (and names)
- ▶ What is the problem and data you worked on?
- ▶ What were the key ideas in your approach?
- ▶ What techniques from the class did you use?
- ▶ What did you learn?

Big Data

Algo data too big for
our computer

(Stats: the limit is here)

(1) Data enormous
200 - 1000 machines

(2) Data Static (some appends)



Algo count

Data Science Revolution

Classic Science

- ① Make hypothesis
- ② Gather Data
- ③ Check if data supports hypothesis

← iid, used once

Stats

④ Precompute, gather ^{store} massive Data

① Find statistically significant differences in data

Data Mining

MapReduce: @ Google

- Built inverted index
- Run PageRank

Distributed File System

<key, value> pair

ISOX

key: "unique" id

value:

raw, somewhat
messy

data

log id

actual log

web address

html, out-going
links

doc id

list of words
↳ k-grams

word index

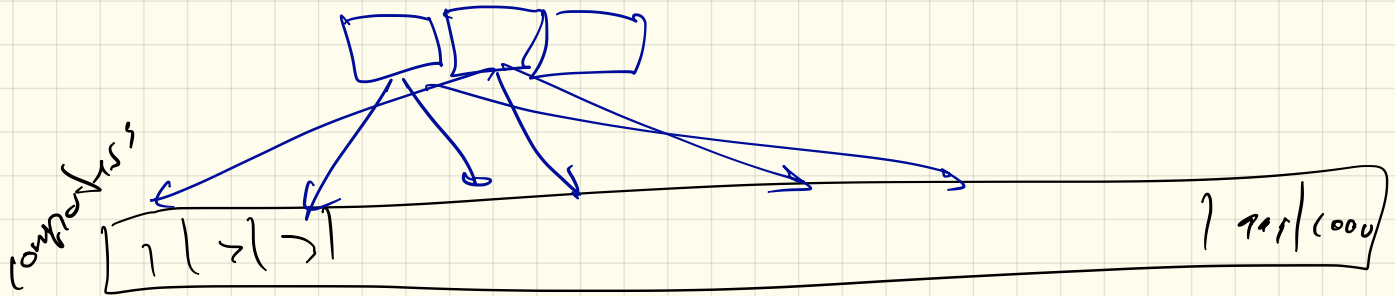
list of docs w/
that word.

Data in Blocks

Block = 64 MB

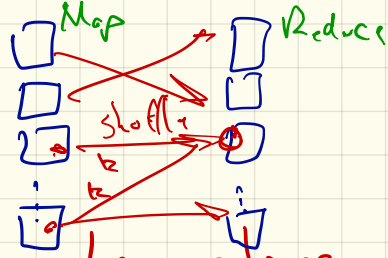
Each block is stored 2 or 3 times

get rid of locality



very resilient (to machine failure)

Map Reduce



1. Map: Read blocks from key-value store

- pull out important info
- determine "new locality" → identify new key

1.5 Combine: Merge $\langle k, v \rangle$ before shuffle

2. Shuffle group by key

3. Reduce: → same key, same new block
look at all $\langle k, v \rangle$ pairs w/ same
new key, post-process
(analysis)

Word Count

(Hello World)

Consider as input all of English Wikipedia stored in DFS. Goal is to count how many times each word is used.

Map: doc id \rightarrow $\{ \langle \overset{it}{word}, id \rangle, \langle \overset{was}{word}, id \rangle, \langle \overset{good}{word}, id \rangle \}$

Combine: $\{ \langle w, v_1 \rangle, \langle w, v_2 \rangle, \dots \} \rightarrow \langle w, \sum_i v_i \rangle$

Reduce: Input $\{ \langle good, v_1 \rangle, \langle good, v_2 \rangle, \dots, \langle good, v_n \rangle \}$
Output $\langle good, \sum_i v_i \rangle$

"the" file

Input $\{ \langle the, 1 \rangle, \langle the, 1 \rangle, \dots, \langle the, 1 \rangle \}$

Inverted Index

Consider as input all of English Wikipedia stored in DFS. Goal is to build an index, so each word has a list of pages it is in.

Map: $\text{article}_{id} \rightarrow \langle \text{word}_1, id \rangle, \langle \text{word}_2, id \rangle, \dots$

Reduce: $\langle \text{word}_i, id_1 \rangle, \langle \text{word}_i, id_2 \rangle, \dots$
 $\rightarrow \langle \text{word}_i, \bigcup_j id_j \rangle$

Phrases

Consider as input all of English Wikipedia stored in DFS. Goal is to build an index, on 3-grams (sequence of 3 words) that appears on exactly one page, with link to page.

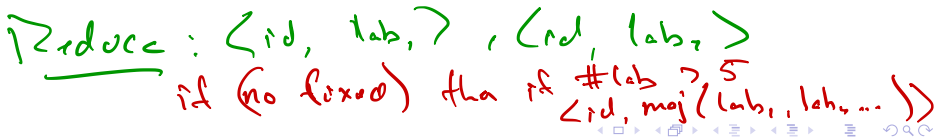
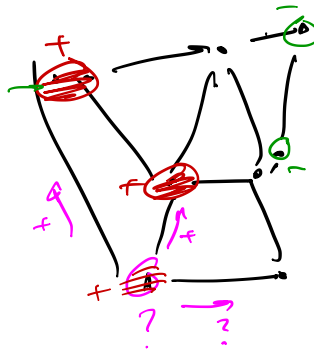
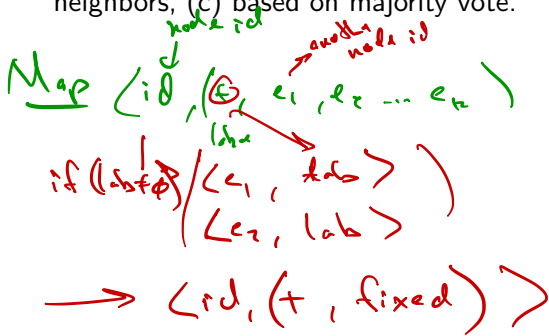
Map : doc : $\langle \text{word}_1, \text{word}_2, \text{word}_3, \text{id} \rangle$
id : $\langle \text{word}_2, \text{word}_3, \text{word}_{21}, \text{id} \rangle$
:

Reduce: $\langle (w_1, w_2, w_{11}), \text{id}_1 \rangle, \langle (w_1, w_2, w_{11}), \text{id}_2 \rangle$
 $\rightarrow \langle (w_1, w_2, w_{11}), \bigcup_i \text{id}_i \rangle$

Label Propagation (Graph)

Consider a large graph $G = (V, E)$ (e.g., a social network), with a subset of nodes $V' \subset V$ with labels (e.g., {pos, neg}). Each node stores its label (if any) and edges.

Assign a vertex a label if (a) unlabeled, (b) has ≥ 5 labeled neighbors, (c) based on majority vote.



Label Propagation (Embedding)

Consider a data set $X \subset \mathbb{R}^d$, with a subset of points $X' \subset X$ with labels (e.g., {pos, neg}). Implicitly defines graph with $V = X$ and E using $k = 20$ nearest neighbors.

Assign a vertex a label if (a) unlabeled, (b) has ≥ 5 labeled neighbors, (c) based on majority vote.

LSH, in 2 rounds

Example PageRank

1 block

$$M = \begin{bmatrix} 0 & 1/2 & 0 & 0 \\ 1/3 & 0 & 1 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{bmatrix}$$

one web page

$$g^* = P^x g_1$$

$$= P(P(P \dots (P g_1)))$$

mat-vec

1 round

≤ 50 rounds

Example PageRank

$$M = \begin{matrix} & \begin{matrix} q(1) & q(2) & q(3) & q(4) \end{matrix} \\ \begin{matrix} \left\langle \text{id}, q(i) \cdot \frac{1}{3} \right\rangle \\ \leftarrow \\ \leftarrow \\ \leftarrow \end{matrix} & \begin{bmatrix} 0 & 1/2 & 0 & 0 \\ 1/3 & 0 & 1 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{bmatrix} \end{matrix}$$

Stripes:

$$M_1 = \begin{bmatrix} 0 \\ 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad M_2 = \begin{bmatrix} 1/2 \\ 0 \\ 0 \\ 1/2 \end{bmatrix} \quad M_3 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad M_4 = \begin{bmatrix} 0 \\ 1/2 \\ 1/2 \\ 0 \end{bmatrix}$$

These are stored as $(1 : (1/3, 2), (1/3, 3), (1/3, 4))$,
 $(2 : (1/2, 1), (1/2, 4))$, $(3 : (1, 2))$, and $(4 : (1/2, 1), (1/2, 3))$.

Example PageRank

$$M = \begin{bmatrix} 0 & 1/2 & 0 & 0 \\ 1/3 & 0 & 1 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{bmatrix}$$

Handwritten annotations: $g(1)$ and $g(2)$ above the first two columns, id_4 with an arrow pointing to the bottom row, and red and green boxes highlighting specific elements and rows.

Blocks:

$$M_{1,1} = \begin{bmatrix} 0 & 1/2 \\ 1/3 & 0 \end{bmatrix} \quad M_{1,2} = \begin{bmatrix} 0 & 0 \\ 1 & 1/2 \end{bmatrix} \quad M_{2,1} = \begin{bmatrix} 1/3 & 0 \\ 1/3 & 1/2 \end{bmatrix} \quad M_{2,2} = \begin{bmatrix} 0 & 1/2 \\ 0 & 0 \end{bmatrix}$$

These are stored as $(1 : (1/2, 2))$, $(2 : (1/3, 1))$, as $(2 : (1, 3), (1/2, 4))$, as $(3 : (1/3, 1))$, $(4 : (1/3, 1), (1/2, 2))$, and as $(3 : (1/2, 4))$.

$$\langle id_4, \overbrace{g(1) \cdot \frac{1}{3} + g(2) \cdot \frac{1}{2}} \rangle$$