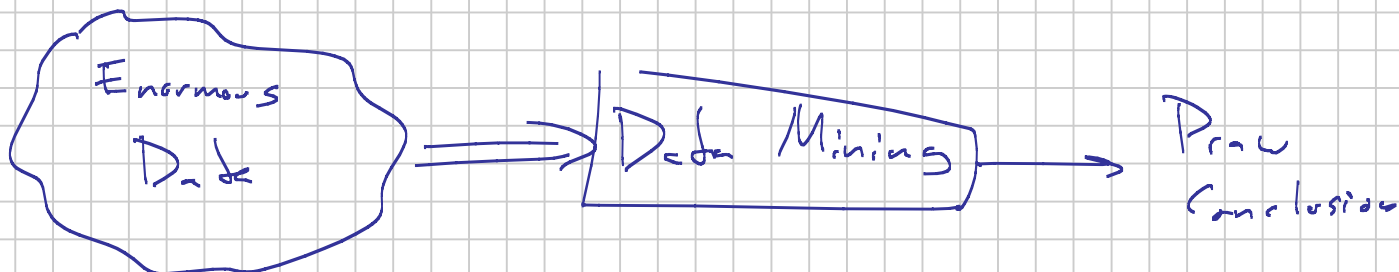
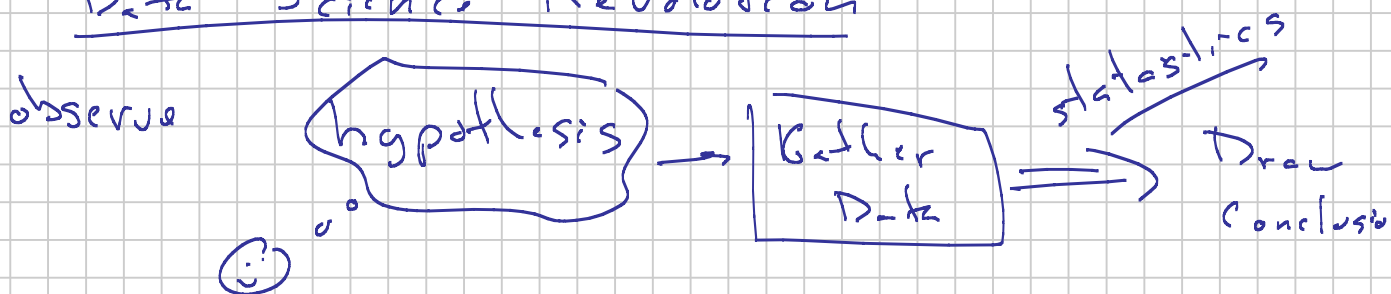


MapReduce + DFS

Big Data

- Data in Enormous ... \geq Terabytes
 \hookrightarrow 100 or 1000s machines
- Data static. Usually, appends.

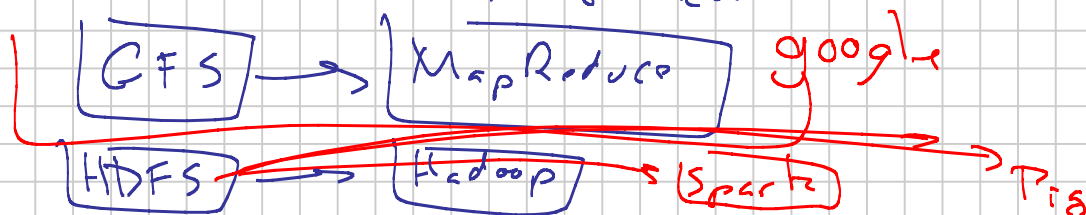
Data Science Revolution

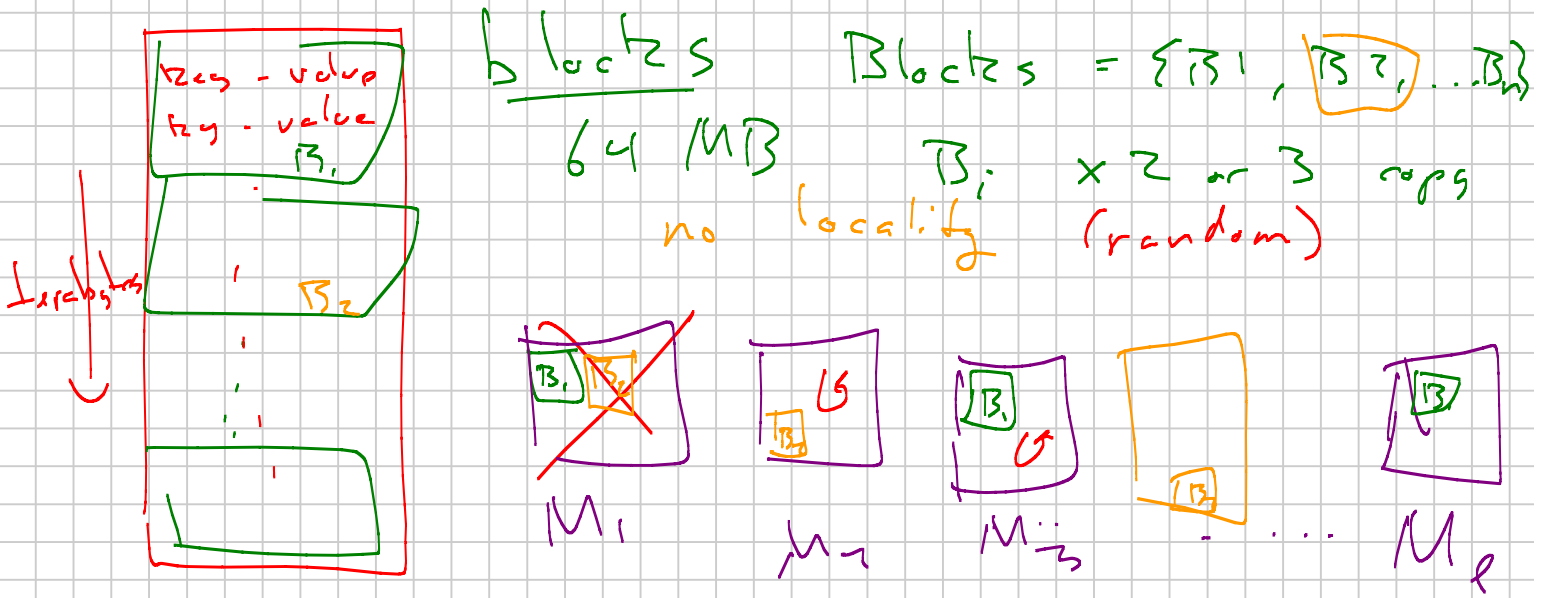


Distributed File Systems

key-value pairs

key	value
log id	log
web address	html, outgoing links
doc id	list of words
word	list of documents containing word





- Redundancy (Disk failure) \rightarrow Resiliency
- Heterogeneity

Map Reduce

1. Map
 - 1.5 Combine (Reduce by shuffle)

Operates on data in DFS
 $\text{map}(\langle k, v \rangle) \rightarrow \{\langle k_i, v_i \rangle, \langle k_i, v_i' \rangle, \dots\}$
2. Shuffle

Put all $\langle k_1, v_1 \rangle, \langle k_2, v_2 \rangle$ s.t.
 $k_1 = k_2$ on same machine
3. Reduce

$\langle k, \{v_1, v_2, v_3, \dots, v_m\} \rangle \rightarrow \langle k, v \rangle$

Word Count

Input: large set of documents = { words }

Output: For each word: how many times

Map $V = \{w_1, w_2, w_3, \dots, w_n\} \rightarrow \{\langle w_1, 1 \rangle, \langle w_1, 1 \rangle, \dots\}$

Reduce $(k = \text{word}) \{v_1 = 1, v_2 = 1, v_3 = 1, \dots, v_n = 1\}$

→ $\langle \text{word} \mid \sum_{i=1}^{top} v_i \rangle$

word = "the" 7% of all words

Combine: $\langle \text{word} \mid \sum_{i=1}^{top} v_i \rangle$

"curse of the last reducer"

Inverted Index

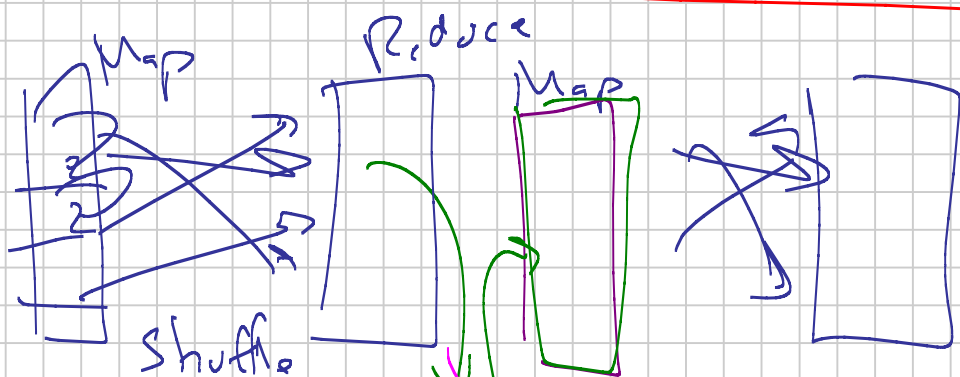
Input: web

Output: $\langle \text{word}, \text{list pages} \rangle$

Map: html → $\left\{ \begin{array}{l} \langle \text{word}_1, (p, \text{info}) \rangle \\ \langle \text{word}_2, (p, \text{info}) \rangle \\ \vdots \\ \langle \text{word}_n, (p, \text{info}) \rangle \end{array} \right\}$

Reduce: $\langle \text{word}_i, (p_1, \text{info}), (p_2, \text{info}), (p_3, \text{info}) \dots \rangle$
→ word & sort (p_1, p_2, \dots, p_n)

Rounds



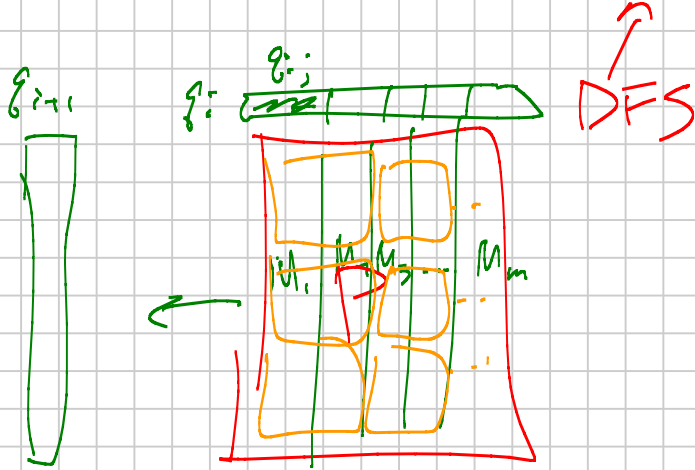
Big overhead

Spark RDD

Page Ranks on MR

So rounds

$$g_{i+1} = \left((1-\beta)P + \beta \frac{1}{n} \right) g_i$$



$$\text{Map: } r_j = M_j g_i$$

$$\text{Reduce: } g_{i+1} = \sum_j r_j (1-\beta) + \frac{\beta}{n}$$