

Implementations to deal with huge data sets.

Note: For huge data sets hardware errors will occur almost certainly.

Map Reduce and Hadoop

(Not the main focus of this lecture, only briefly summarized)

Key idea:

Use parallelism from "computing clusters" (not a super-computer), built of commodity hardware, connected by Ethernet and inexpensive switches.

Software stack:

(i) Distributed file system (DFS)

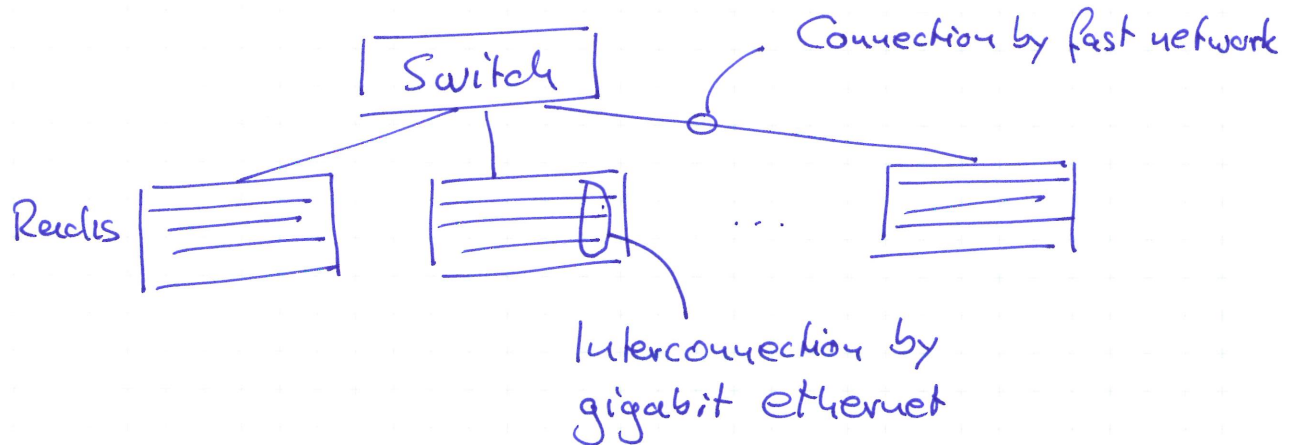
- large blocks
- redundancy by replication

(ii) Programming system: MapReduce

- tolerant to hardware failure
- able to handle large data sets efficiently

Architecture:

- (i) Compute nodes stored on racks, each with its own processor and storage device.
- (ii) Racks are connected by switches



Principles -

- (i) Files are stored redundantly to protect against failure of nodes.
- (ii) Computations are divided into independent tasks. If one fails it can be restarted without affecting others.

Remarks (i): Distributed File system (DFS)

- o Files are divided into chunks (typically 64MB)
- o Chunks are replicated (typically 3 times on different racks)
- o A file master node or name node has information where to find copies of files

Implementations :

- o GFS (Google file system)
- o HDFS (Hadoop distributed file system, Apache)
- o Cloud Store (open source DFS)

Remarks (ii) Map Reduce (computing paradigm)

- o System manages parallel execution and coordination of tasks.
- o 2 functions are written by the user: Map and Reduce

Implementations :

- o MapReduce (Google, internal)
- o Hadoop (Open source, Apache)

<http://hadoop.apache.org/>

What Is Apache Hadoop?

The Apache™ Hadoop® project develops open-source software for reliable, scalable, distributed computing.

The Apache Hadoop software library is a framework that allows for the **distributed processing of large data sets across clusters of computers** using simple programming models. It is designed to **scale up from single servers to thousands of machines**, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to **detect and handle failures at the application layer**, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

The project includes these modules:

- **Hadoop Common:** The common utilities that support the other Hadoop modules.
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.
- **Hadoop YARN:** A framework for job scheduling and cluster resource management.
- **Hadoop MapReduce:** A YARN-based system for parallel processing of large data sets.

Co-founders: Doug Cutting and Mike Cafarella, January 2006

(Doug Cutting named the system after his son's toy elephant.)

2. Prerequisites from Matrix Algebra

Real $(n \times n)$ matrices will be written as

$$M = (m_{ij})_{\substack{1 \leq i \leq n \\ 1 \leq j \leq n}} \in \mathbb{R}^{n \times n} \quad (\text{or } \mathbb{C}^{n \times n})$$

Diagonal matrices as $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \in \mathbb{R}^{n \times n}$

Matrix $U \in \mathbb{R}^{n \times n}$ is called orthogonal if

$$UU^T = U^T U = I_n \quad (\text{identity matrix})$$

$O(n)$ set of orthogonal $(n \times n)$ -matrices.

Th. 2.1. (Singular value decomposition, SVD)

Given $M \in \mathbb{R}^{n \times n}$. There exist $U \in O(n)$

and $V \in O(n)$ and some $\Sigma \in \mathbb{R}^{n \times n}$ with non-negative entries in its diagonal and zeros otherwise such that

$$M = U \Sigma V^T.$$

The diagonal elements of Σ are called singular values. The columns of U and V are called left and right singular vectors.

+

Remark If $m < n$

$$\S \quad \boxed{M} = \boxed{U} \boxed{\begin{matrix} \circ & & & \\ & \circ & & \\ & & \circ & \\ & & & \circ \end{matrix}} \boxed{V^T}$$

then the SVD may be written as

$$\exists U \in \mathbb{R}^{m \times m}, U U^T = I_m, \exists V \in O(n), \exists \Sigma \in \mathbb{R}^{m \times m}$$

diagonal with ≥ 0 entries

such that $M = U \Sigma V$. \perp

Th 2.2. (Spectral decomposition)

Given $M \in \mathbb{R}^{n \times n}$ symmetric. There exist $V \in O(n)$,
 $V = (v_1, \dots, v_n)$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ such that

$$M = V \Lambda V^T = \sum_{i=1}^n \lambda_i v_i v_i^T.$$

v_i are the eigenvectors of M with eigenvalues λ_i . \perp

o If $\lambda_i > 0, i=1, \dots, n$, M is called positive definite
 $(M > 0)$ (p.d.)

If $\lambda_i \geq 0, i=1, \dots, n$, M is called non-negative definite
 $(M \geq 0)$ (n.n.d.)

o If $M \geq 0$, then it has a Cholesky decomposition

$$M = V \Lambda^{\frac{1}{2}} (V \Lambda^{\frac{1}{2}})^T$$

where $\Lambda^{\frac{1}{2}} = \text{diag}(\lambda_1^{\frac{1}{2}}, \dots, \lambda_n^{\frac{1}{2}})$

o $M \succeq 0 \Leftrightarrow x^T M x \geq 0 \quad \forall x \in \mathbb{R}^n$

$$\left[\begin{array}{l} \Rightarrow M \succeq 0 \Rightarrow M = V \Lambda V^T, \Lambda \succeq 0 \\ \Rightarrow x^T M x = x^T V \Lambda V^T x \geq 0 \quad \forall x \in \mathbb{R}^n. \end{array} \right]$$

o $M \succ 0 \Leftrightarrow x^T M x > 0 \quad \forall x \in \mathbb{R}^n, x \neq 0.$

Def. 2.3. a) $M = (m_{ij}) \in \mathbb{R}^{n \times n}$. $\text{tr}(M) = \sum_{i=1}^n m_{ii}$ is called trace of M .

b) Given $M \in \mathbb{R}^{n \times n}$. $\|M\|_F = \left(\sum_{i,j} m_{ij}^2 \right)^{1/2} = \left(\text{tr}(M^T M) \right)^{1/2}$ is called the Frobenius norm of M .

c) $M \in \mathbb{R}^{n \times n}$, M symmetric. $\|M\|_S = \max_{1 \leq i \leq n} |\lambda_i|$ is called spectral norm.

o It holds that $\text{tr}(A \cdot B) = \text{tr}(B \cdot A)$, $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times n}$.

o $\text{tr}(M) = \sum_{i=1}^n \lambda_i(M)$, $\det(M) = \prod_{i=1}^n \lambda_i(M)$ for any symm. matrix M .

Th. 2.4. (Ky Fan, 1950)

Given $M \in \mathbb{R}^{n \times n}$ symm., $k \leq n$, $\lambda_1(M) \geq \dots \geq \lambda_n(M)$ eigenvalues.

$$\max_{\substack{V \in \mathbb{R}^{n \times k} \\ V^T V = I_k}} \text{tr}(V^T M V) = \sum_{i=1}^k \lambda_i(M)$$

$$\min_{\substack{V \in \mathbb{R}^{n \times k} \\ V^T V = I_k}} \text{tr}(V^T M V) = \sum_{i=1}^k \lambda_{n-i+1}(M)$$

Special case: $k=1$

$$\max_{\|v\|=1} v^T M v = \lambda_{\max}(M)$$

$$\min_{\|v\|=1} v^T M v = \lambda_{\min}(M)$$

Also note:

$$\max_{\|v\|=1} v^T M v = \max_{v \neq 0} \frac{v^T M v}{v^T v}$$

Th. 2.5. Given $A, B \in \mathbb{R}^{n \times n}$, symmetric with eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ and $\mu_1 \geq \dots \geq \mu_n$, respectively. Then

$$\sum_{i=1}^n \lambda_i \mu_{n-i+1} \leq \text{tr}(A \cdot B) \leq \sum_{i=1}^n \lambda_i \mu_i$$

Proof. First show for any $x_1, \dots, x_n, y_1, \dots, y_n \in \mathbb{R}$
with ordered values $x_{[1]} \geq \dots \geq x_{[n]}$ and $y_{[1]} \geq \dots \geq y_{[n]}$.

$$\sum_{i=1}^n x_{[i]} y_{[n-i+1]} \leq \sum_{i=1}^n x_i y_i \leq \sum_{i=1}^n x_{[i]} y_{[i]} \quad (*)$$

Left as an exercise. Now:

It holds

$$\begin{aligned} \text{tr}(A \cdot B) &= \text{tr}(V \Lambda V^T U M U^T) \\ &= \text{tr}(\underbrace{U^T V}_Q \Lambda \underbrace{V^T U}_Q M) \quad , \quad Q = (q_1, \dots, q_n) = U^T V \\ &\quad \text{orthogonal} \\ &= \text{tr}((\lambda_1 q_1, \dots, \lambda_n q_n) \begin{pmatrix} \mu_1 q_1^T \\ \vdots \\ \mu_n q_n^T \end{pmatrix}) \\ &= \text{tr}(\sum_{i=1}^n \lambda_i \mu_i q_i q_i^T) \\ &= \sum_{i=1}^n \lambda_i \mu_i \text{tr}(q_i q_i^T) \\ &= \sum_{i=1}^n \lambda_i \mu_i \underbrace{\text{tr}(q_i^T q_i)}_{=1} = \sum_{i=1}^n \lambda_i \mu_i . \end{aligned}$$

Assertion follows by (*). \square

Let $\lambda^+ = \max\{\lambda, 0\}$ denote the positive part of $\lambda \in \mathbb{R}$.

Th. 2.6. Given $M \in \mathbb{R}^{n \times n}$ symmetric with spectral decomposition $M = V \text{diag}(\lambda_1, \dots, \lambda_n) V^T$, $\lambda_1 \geq \dots \geq \lambda_n$. Then for $k \leq n$

$$\min_{A \geq 0, \text{rk}(A) \leq k} \|M - A\|_F^2$$

is attained at $A^* = V \text{diag}(\lambda_1^+, \dots, \lambda_k^+, 0, \dots, 0) V^T$ with optimum value $\sum_{i=1}^k (\lambda_i - \lambda_i^+)^2 + \sum_{i=k+1}^n \lambda_i^2$. \perp

Proof.

$$\|M - A\|^2 = \|M\|^2 - 2\text{tr}(MA) + \|A\|^2$$

$$\geq \sum_{i=1}^n \lambda_i^2 - 2 \sum_{i=1}^n \lambda_i \mu_i + \sum_{i=1}^n \mu_i^2$$

$\mu_1 \geq \dots \geq \mu_n \geq 0$ eigenv. of A

$$= \sum_{i=1}^n (\lambda_i - \mu_i)^2$$

$$= \sum_{i=1}^k (\lambda_i - \mu_i)^2 + \sum_{i=k+1}^n (\lambda_i - 0)^2, \text{ if } \text{rk}(A) \leq k$$

$$\geq \sum_{i=1}^k (\lambda_i - \lambda_i^+)^2 + \sum_{i=k+1}^n \lambda_i^2$$

Lower bound is attained if $A = V \text{diag}(\lambda_1^+, \dots, \lambda_k^+, 0, \dots, 0) V^T$. \square