

Lecture Notes 1: Matrix multiplication

Lecturer: Che-Rung Lee

Scribe: Tien-Yu Kuo 9962526

1 Matrix multiplication

1.1 Basic definition

- Direct multiplication

$$\text{Let } A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mp} \end{bmatrix}, B = \begin{bmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ b_{21} & b_{22} & \cdots & b_{2n} \\ \vdots & \vdots & & \vdots \\ b_{p1} & b_{p2} & \cdots & b_{pn} \end{bmatrix} \text{ where } \begin{cases} \dim(A) = (m, p) \\ \dim(B) = (p, n) \end{cases}$$

$$\text{Then } C = A \cdot B = \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1n} \\ c_{21} & c_{22} & \cdots & c_{2n} \\ \vdots & \vdots & & \vdots \\ c_{m1} & c_{m2} & \cdots & c_{mn} \end{bmatrix}$$

$$\text{where } c_{ij} = \sum_{k=1}^p a_{ik}b_{kj}, \forall i = 1, \dots, m, j = 1, \dots, n, \dim(C) = (m, n) \quad (1)$$

$$\text{Complexity : } O(mnp) \begin{cases} \text{multiplication \# : } mnp \\ \text{addition \# : } mn(p-1) \end{cases}$$

- Inner Product form

Generally, we treat all vectors as column vectors.

Let \vec{a}, \vec{b} with $\dim(\vec{a}) = \dim(\vec{b}) = n$

$$\Rightarrow \vec{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}, \vec{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \text{ and the inner product } (\vec{a} \cdot \vec{b}) = a_1b_1 + \dots + a_nb_n = \vec{a}^T \vec{b}$$

$$\text{Hence, if } A = \begin{bmatrix} a_{11} & \cdots & a_{1p} \\ a_{21} & \cdots & a_{2p} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mp} \end{bmatrix} = \begin{bmatrix} \vec{a}_1^T \\ \vec{a}_2^T \\ \vdots \\ \vec{a}_m^T \end{bmatrix}, B = \begin{bmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ \vdots & \vdots & & \vdots \\ b_{p1} & b_{p2} & \cdots & b_{pn} \end{bmatrix} = [\vec{b}_1 \quad \vec{b}_2 \quad \cdots \quad \vec{b}_n]$$

$$\text{Then } C = A \cdot B = \begin{bmatrix} \vec{a}_1^T \vec{b}_1 & \vec{a}_1^T \vec{b}_2 & \cdots & \vec{a}_1^T \vec{b}_n \\ \vdots & \vdots & & \vdots \\ \vec{a}_m^T \vec{b}_1 & \vec{a}_m^T \vec{b}_2 & \cdots & \vec{a}_m^T \vec{b}_n \end{bmatrix}$$

$$\text{where } c_{ij} = \vec{a}_i^T \vec{b}_j, i = 1, \dots, m, j = 1, \dots, n \quad (2)$$

- Outer Product form

$$\text{Let } A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mp} \end{bmatrix} = [\vec{a}_1 \quad \vec{a}_2 \quad \cdots \quad \vec{a}_p], B = \begin{bmatrix} b_{11} & \cdots & b_{1n} \\ b_{21} & \cdots & b_{2n} \\ \vdots & & \vdots \\ b_{p1} & \cdots & b_{pn} \end{bmatrix} = \begin{bmatrix} \vec{b}_1^T \\ \vec{b}_2^T \\ \vdots \\ \vec{b}_p^T \end{bmatrix}$$

$$\text{Then } C = A \cdot B = \sum_{k=1}^p \vec{a}_k \vec{b}_k^T.$$

The elements in C are the same as (1).

Proof:

$$\text{Let } C = C^{(1)} + \cdots + C^{(k)} + \cdots + C^{(p)}, C^{(k)} = \vec{a}_k \vec{b}_k^T$$

$$\therefore c_{ij}^{(k)} = a_{ik} b_{kj} \Rightarrow c_{ij} = \sum_{k=1}^p c_{ij}^{(k)} = \sum_{k=1}^p a_{ik} b_{kj}$$

...the same as (1)

$$\text{e.g. } A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, B = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \Rightarrow C = C^{(1)} + C^{(2)} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \begin{bmatrix} 4 \\ 8 \end{bmatrix} = \begin{bmatrix} 5 \\ 11 \end{bmatrix}$$

Review Rank of a matrix.

Column(row) rank : the maximum number of linearly independent column(row) vectors.

The column rank and the row rank are always equal and hence it is simply called the rank of a matrix.

- Block form

$$A = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1P} \\ \vdots & \vdots & & \vdots \\ A_{M1} & A_{M2} & \cdots & A_{MP} \end{bmatrix} \text{ which contains } P \text{ (} M \text{) blocks per row (column).}$$

e.g. If $\dim(A) = (1024, 1024)$, $\dim(A_{ij})$ could be $(64, 64) \cdots$ all the same

$$\text{e.g. } \begin{bmatrix} \dim(5, 6) & \dim(5, 8) & \dim(5, 7) \\ \dim(3, 6) & \dim(3, 8) & \dim(3, 7) \\ \dim(2, 6) & \dim(2, 8) & \dim(2, 7) \end{bmatrix} \text{ whose all boxes in the same row(column) have}$$

the same height(width).

$$\text{Let } B = \begin{bmatrix} B_{11} & \cdots & B_{1N} \\ \vdots & & \vdots \\ B_{P1} & \cdots & B_{PN} \end{bmatrix}, \text{ if we want to do block form multiplication, } \dim(A_{IK})_2 =$$

$\dim(B_{KJ})_1$ should be true $\forall 1 \leq I \leq M, 1 \leq J \leq N, 1 \leq K \leq P.$

That is, $C = A \cdot B = \begin{bmatrix} C_{11} & C_{12} & \cdots & C_{1N} \\ \vdots & \vdots & & \vdots \\ C_{M1} & C_{M2} & \cdots & C_{MN} \end{bmatrix}, C_{IJ} = \sum_{k=1}^P A_{IK} B_{KJ}$

Question How to prove it the same as (2)? (Homework problem)

1.2 Performance Consideration

- Performance
 - (1) computation
 - (2) bandwidth
 - (3) memory latency

Computation $\rightarrow T_c$, bandwidth, memory latency $\rightarrow T_d$.

Assume two $n \times n$ matrices A and B , then $C = A \cdot B$ is also $n \times n$. \Rightarrow Totally $3n^2$ data to be stored.

$\therefore T_d$ is related to $3n^2$, T_c is related to $2n^3$.

What if the storing spaces are much less?

If we only have $3b^2$ spaces, $b \ll n$, in our "fast" memory and assume $b^2 \approx n$. In order to compute C , A will be loaded n times, each of them is a complete A (i.e. n^2 elements), as shown in Fig. 1. So the total data movement is $O(n^3)$.

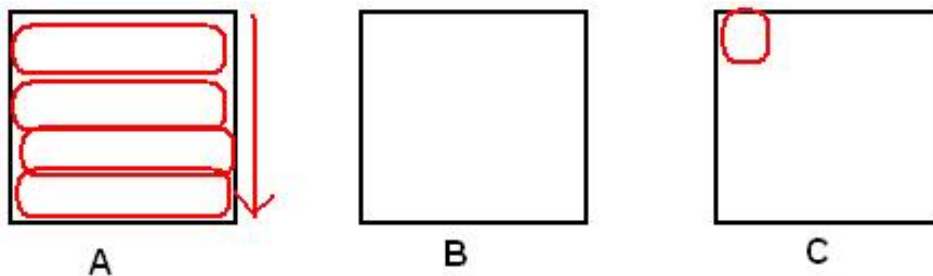


Figure 1: If the spaces are much less.

How about using block form?

Computing a block C_{ij} needs $2b^2 \times \frac{n}{b} = 2bn$ data movement and $2b^3 \times \frac{n}{b} = 2b^2n$ computations, as shown in Fig. 2.

\therefore We have $(\frac{n}{b})^2 C_{ij}$ blocks and hence require $2bn \times \frac{n^2}{b^2} = \frac{2n^3}{b}$ loadings, $2b^2n \times \frac{n^2}{b^2} = 2n^3$ computations. $\Rightarrow T_c = 2n^3, T_d = \frac{2n^3}{b} \dots$ b times less data movement.

- A brief history of (Dense) Linear Algebra software
 - BLAS (1973-1977)
 - Also called BLAS 1, do $O(n^1)$ ops on $O(n^1)$ data. Standard library of 15 operations (mostly) on vectors.

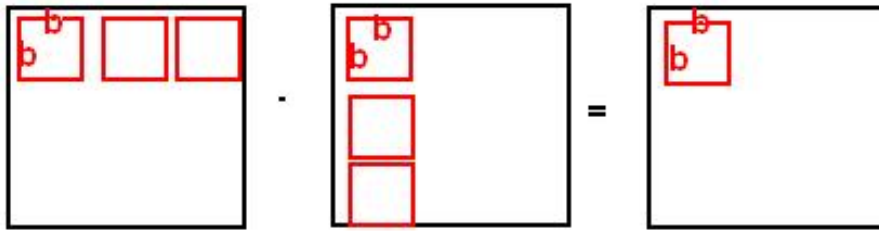


Figure 2: Using block form.

- BLAS 2 (1984-1986)
Standard library of 25 operations (mostly) on matrix/vector pairs. Do $O(n^2)$ ops on $O(n^2)$ data.
- BLAS 3 (1987-1988)
Standard library of 9 operations (mostly) on matrix/matrix pairs, functions are integrated. Do $O(n^3)$ ops on $O(n^2)$ data...much higher computational intensity.
- LAPACK (1989-now)
"Linear Algebra PACKage"-uses BLAS-3. Much more generalized.
- ScaLAPACK (1995-now)
"Scalable LAPACK", uses MPI to parallelize. More complex data structures, algorithms than LAPACK.