

VE472 Lecture 13

Jing Liu

UM-SJTU Joint Institute

Summer

- Many methods can be understood using various matrix decompositions, thus approximating the product of two matrices is a way to deal with large data.

$$\mathbf{AB} \approx \mathbf{GF}$$

- Intuitively, it means \mathbf{AB} is close to \mathbf{GF} , but we need to define it precisely.
- Recall the notion of norm allows to talk about $\mathbf{x} \in \mathbb{R}^n$ is close to $\mathbf{y} \in \mathbb{R}^n$

$$\mathbf{x} \approx \mathbf{y} \quad \text{if} \quad \|\mathbf{x} - \mathbf{y}\| < \varepsilon$$

for some small positive $\varepsilon \in \mathbb{R}$, and the norm of $\mathbf{x} - \mathbf{y}$ give how close they are.

- Therefore, we need to discuss matrix norms as well as vector norms!
- We have used largely the ℓ_2 -norm for any $\mathbf{u} \in \mathbb{R}^n$ in this course,

$$\|\mathbf{u}\| = \sqrt{\sum_{i=1}^n u_i^2} \quad \text{vector norm}$$

- For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, one useful matrix norm is the Frobenious norm

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$$

- Another important matrix norm is given by

operator norm $\|\mathbf{A}\|_O = \max_{\hat{\mathbf{x}}} \|\mathbf{A}\hat{\mathbf{x}}\|$ where $\|\hat{\mathbf{x}}\| = 1$ ^{unit vector}
vector norm

which is called the operator norm induced by the vector norm $\|\cdot\|$ on \mathbb{R}^n .

- It can be shown the two matrix norms are related to each other as

$$\|\mathbf{A}\|_O \geq \|\mathbf{A}\|_F \geq \sqrt{n} \|\mathbf{A}\|_O$$

- So the problem we want to solve is whether we can find a \mathbf{G} and \mathbf{F} such that

$$\|\mathbf{AB} - \mathbf{GF}\|_F \quad \text{or} \quad \|\mathbf{AB} - \mathbf{GF}\|_O$$

is small, and computationally cheaper to find and compute \mathbf{GF} in general.

Algorithm 1: Naive three-loop algorithm for matrix multiplication

Input : Matrix **A** of $m \times p$, and matrix **B** of $p \times n$.

Output: The product $\mathbf{M} = \mathbf{AB}$

1 **Function** NaiveMatrixMultiplication(**A**, **B**):

```
2   for  $i \leftarrow 1$  to  $m$  do      rows of first matrix
3       for  $j \leftarrow 1$  to  $n$  do  columns of second matrix
4            $[\mathbf{M}]_{ij} \leftarrow 0$  ;
5           for  $\ell \leftarrow 1$  to  $p$  do
6                $[\mathbf{M}]_{ij} \leftarrow [\mathbf{M}]_{ij} + [\mathbf{A}]_{i\ell}[\mathbf{B}]_{\ell j}$  ;
7           end for
8       end for
9   end for
10  return  $\mathbf{M}$  ;
11 end
```

- Notice the naive approach essentially computes the matrix multiplication as

$$m \times n \times p$$

pairs of products that need to be summed together, or as

$$m \times n$$

pairs of inner/dot products involving row and column vectors in \mathbb{R}^p .

- The idea is to take a sample of all computations needed to reduce the cost!
- Some of those computations are more important than other others so should be selected with a higher probability, we could also gain the insight from

$$A = \sum_{i=1}^n x_i = \mu \times n \approx \hat{x} \times L$$

← how many samples taken

where μ is the population mean while \hat{x} is the sample mean.

- We could sample according to any probability distribution, including simple random sampling, and obtain an unbiased estimator of the mean/sum, but...
data大的时候可以让bias大一点->更小的var

- In order to better implement this idea of using a non-uniform distribution to minimise the variance, we need to see the matrix multiplication as

see notes

$$\overset{p \times n}{\mathbf{A}} \mathbf{B} = \sum_{k=1}^p \mathbf{a}_k B_k$$

matrix multiplication by outer product

where \mathbf{a}_k denotes the k th column of \mathbf{A} , while B_k denotes the k th row of \mathbf{B} .

- This essentially decomposes the product into p outer products, each needs

$$m \times n$$

scalar multiplications, the idea is to find a sample of those p outer products.

$$\mathbf{A}\mathbf{B} = \sum_{k=1}^p \mathbf{a}_k B_k \approx \frac{1}{L} \sum_{\ell=1}^L \frac{1}{q_{k_\ell}} \mathbf{a}_{k_\ell} B_{k_\ell} = \sum_{\ell=1}^L \mathbf{g}_{k_\ell} F_{k_\ell} = \mathbf{G}\mathbf{F}$$

sample size probability of a particular block being calculated

Q: What are the mean and variance of the ij th element of $[\mathbf{G}\mathbf{F}]$?

Algorithm 2: Basic sampling algorithm for matrix multiplication

Input : Matrix \mathbf{A} of $m \times p$, matrix \mathbf{B} of $p \times n$, positive integer L , and list of probabilities $\{q_i\}_{i=1}^p$.

Output : Matrices \mathbf{G} and \mathbf{F} such that $\mathbf{GF} \approx \mathbf{AB}$

```
1 Function BasicSamplingMatrixMultiplication( $\mathbf{A}, \mathbf{B}$ ):
2   for  $\ell \leftarrow 1$  to  $L$  do
3     Sample  $k_\ell \in \{1, \dots, p\}$  with probability  $\Pr(k_\ell = i) = q_i$ 
4     /* The above step is done independently and identically
5        with replacement from one  $\ell$  value to another. */
6     for  $i \leftarrow 1$  to  $m$  do
7       |  $[\mathbf{G}]_{i\ell} \leftarrow [\mathbf{A}]_{i\ell} / \sqrt{Lq_{k_\ell}}$  ;
8     end for
9     for  $j \leftarrow 1$  to  $n$  do
10      |  $[\mathbf{F}]_{\ell j} \leftarrow [\mathbf{B}]_{\ell j} / \sqrt{Lq_{k_\ell}}$  ;
11    end for
12  end for
13  return  $\mathbf{M}$  ;
14 end
```

Theorem 0.1

Suppose matrices \mathbf{A} and \mathbf{B} , and the matrices \mathbf{G} and \mathbf{F} generated by Algo 2, then

$$\mathbb{E} [\|\mathbf{AB} - \mathbf{GF}\|_F^2] = \sum_{k=1}^p \frac{\|\mathbf{a}_k\|^2 \|B_k\|^2}{Lq_k} - \frac{1}{L} \|\mathbf{AB}\|_F^2$$

expected error

Furthermore, if

$$q_k = \frac{\|\mathbf{a}_k\| \|B_k\|}{A} \quad \text{where} \quad A = \sum_{k=1}^p \|\mathbf{a}_k\| \|B_k\|$$

then the expected Frobenious norm of the error is minimised, and the value is

$$\mathbb{E} [\|\mathbf{AB} - \mathbf{GF}\|_F^2] = \frac{1}{L} \left(\sum_{k=1}^p \|\mathbf{a}_k\| \|B_k\| \right)^2 - \frac{1}{L} \|\mathbf{AB}\|_F^2$$

- Notice the expected error is the sum of all variances $\text{Var} [[\mathbf{GF}]_{ij}]$.
error of using GF instead of AB is sum of GF's variance