

# Low Rank Adaptation(LoRA)

## WHAT?

Approximate representation of a matrix

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

Original matrix

$\approx$

$$\begin{bmatrix} 2.99 & 2.01 & 0.98 \\ 0.02 & 1.88 & 1.1 \\ -0.07 & 0.45 & 0.29 \end{bmatrix}$$

The numbers in this matrix are very close to the original matrix

# Low Rank Adaptation (LoRA)

## WHY?

Why do we want to approximate a matrix?

Bigger matrix → Lots of numbers (or elements) → need more storage space and more computation

*How about finding a way to store the matrix with less number of elements with approximate same values?*

Yes! There is a way

We can store the matrix with much less elements → Much less storage space & Much less computation power needed

# Low Rank Adaptation(LoRA)

## How?

Through “**Rank of a Matrix**”

Rank= No. of independent columns or rows of a matrix

Let A be a matrix of order (m x n) and its rank is 'r'. [ $r \leq \min(m,n)$ ]

We can split the matrix A into two matrices through 'r' or even any value ranging from 1 to r.

$$A_{m \times n} = B_{m \times r} C^T_{r \times n}$$

Example

Let m=No. of rows=5000

n=No. of columns=100

r=rank of A=20

( $A_{5000 \times 100}$ )

Now Number of elements of A =  $5000 \times 100 = 5,00,000$

We can find two matrices B and  $C^T$  such that

$B_{5000 \times 20}$  and  $C^T_{20 \times 100}$

No. of elements of  $BC^T = 5000 \times 20 + 20 \times 100 = 1,02,000$

$$A_{5000 \times 100} = B_{5000 \times 20} C^T_{20 \times 100}$$

So instead of storing the original matrix A with 5,00,000 elements we can store B and  $C^T$  matrix with 1,02,000 elements i.e. almost **80%** reduction in storage space required !

# Low Rank Adaptation(LoRA)

## Want further reduction in the storage space requirement?

Choose any lower value between 1 and 20 (i.e. lower rank of A , Full rank is 20)

Lets choose a low rank  $r=10$

Now we can write  $A_{5000 \times 100} = B_{5000 \times 10} C^T_{10 \times 100}$

No. of elements =  $5000 \times 10 + 10 \times 100 = 51,000$

So instead of storing the original matrix A with 5,00,000 elements we can store B and  $C^T$  matrix with 51,000 elements i.e. almost 90% reduction in storage space required !

- This is called Low rank approximation of A, because we are using a low rank of 10, instead of its full rank of 20
- So we can represent the same matrix A with much less number of elements which is very close to the original matrix

# Low Rank Adaptation(LoRA)

**Q:How to find B and C<sup>T</sup> ?**

**A: By Singular Value Decomposition(SVD)**

Any matrix of order (mxn) can be written as a matrix multiplication of 3 matrices as follows:

$A=U\Sigma V^T$  and rank of A is 'r'

Where:

U is a mxm matrix

V is nxn matrix

$\Sigma$  is a mxn matrix, the diagonal values of first 'r' rows are singular values of A and rest of the entries are zero.

Handwritten diagram illustrating the SVD decomposition of a matrix  $A$  into  $U$ ,  $\Sigma$ , and  $V^T$ .

The diagram shows the equation:

$$A = U \Sigma V^T$$

where:

- $A$  is a matrix of order  $m \times n$ .
- $U$  is a matrix of order  $m \times m$ .
- $\Sigma$  is a matrix of order  $m \times n$  with diagonal values  $\sigma_1, \sigma_2, \dots, \sigma_r$ .
- $V^T$  is a matrix of order  $n \times n$ .

Below the diagram, it is noted that:

$$r \rightarrow \text{Rank}(A) \leq \min(m, n)$$

Handwritten example illustrating the SVD decomposition of a matrix  $A$  into  $U$ ,  $\Sigma$ , and  $V^T$ .

The diagram shows the equation:

$$A = U \Sigma V^T$$

where:

- $A$  is a matrix of order  $20 \times 10$ .
- $U$  is a matrix of order  $20 \times 20$ .
- $\Sigma$  is a matrix of order  $20 \times 10$  with diagonal values  $0.8, 0.5, 0.2$ .
- $V^T$  is a matrix of order  $10 \times 10$ .

Additional notes:

- $m = 20 \rightarrow \text{rows}$
- $n = 10 \rightarrow \text{columns}$
- $\text{Rank}(A) = r = 8$
- The 8th row of  $\Sigma$  is highlighted.

Note: SVD can be done by any computer programming language such as python

# Low Rank Adaptation(LoRA)

## Example

$$\text{Let } A = \begin{bmatrix} 3 & 2 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

Here rank of A is 3 as there are three independent columns.  
Let's do a low rank approximation say 2 rank approximation

Applying SVD(Here I have used python to find out U,  $\Sigma$  and  $V^T$  matrices), A can be split as follows:

$$\begin{bmatrix} 3 & 2 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 1 \end{bmatrix} = \underbrace{\begin{bmatrix} 0.91 & 0.42 & 0.02 \\ 0.41 & -0.87 & -0.26 \\ 0.09 & -0.24 & 0.97 \end{bmatrix}}_U \underbrace{\begin{bmatrix} 4.04 & 0 & 0 \\ 0 & 1.70 & 0 \\ 0 & 0 & 0.87 \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} 0.67 & 0.73 & 0.08 \\ 0.65 & -0.54 & -0.53 \\ 0.35 & -0.41 & 0.84 \end{bmatrix}^T}_{V^T}$$

# LowRankAdaptation(LoRA)

## Example

Rank 2 approximation of A:

$$\begin{bmatrix} 0.91 & 0.42 & 0.02 \\ 0.41 & -0.87 & -0.26 \\ 0.09 & -0.24 & 0.97 \end{bmatrix} \begin{bmatrix} 4.04 & 0 & 0 \\ 0 & 1.70 & 0 \\ 0 & 0 & 0.87 \end{bmatrix} \begin{bmatrix} 0.67 & 0.73 & 0.08 \\ 0.65 & -0.54 & -0.53 \\ 0.35 & -0.41 & 0.84 \end{bmatrix}^T$$

$3 \times 2 \quad 2 \times 2 \quad 2 \times 3$

$$= \underbrace{\begin{bmatrix} 0.91 & 0.42 \\ 0.41 & -0.87 \\ 0.09 & -0.24 \end{bmatrix}}_B \underbrace{\begin{bmatrix} 4.04 & 0 \\ 0 & 1.70 \end{bmatrix}}_C \underbrace{\begin{bmatrix} 0.67 & 0.73 & 0.08 \\ 0.65 & -0.54 & -0.53 \end{bmatrix}^T}_{C^T}$$

$$= \begin{bmatrix} 2.99 & 2.01 & 0.98 \\ 0.02 & 1.88 & 1.1 \\ -0.07 & 0.45 & 0.29 \end{bmatrix} \approx \begin{bmatrix} 3 & 2 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 1 \end{bmatrix} \text{ Original Matrix}$$