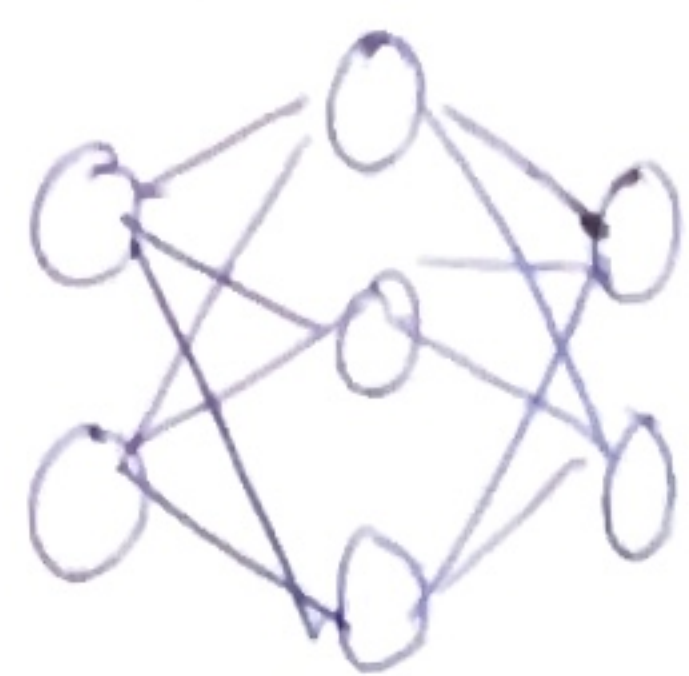


LORA (Type of - Parameter Efficient Fine Tuning)

trains the model by adding more trainable parameters.



$x \rightarrow h(x)$ → hidden layer.

$$h(x) = \underbrace{W_0 x}_{\text{inputs.}}$$

Trainable parameters in a weight matrix (original)

Now say.

$$h(x) = \underbrace{W_0 x}_{\text{keep these frozen (no retraining.)}} + \underbrace{\Delta W x}_{\Delta W = BA \rightarrow \text{new trainable parameter matrix.}}$$

keep these frozen
(no retraining.)

$\Delta W = BA$ → the no. of these parameters are less than W_0 because of low rank decomposition.

$$\text{Params in } B + \text{Params in } A =$$

$$(d \times r) + (r \times k)$$

$$\Rightarrow \textcircled{r} (d + k)$$

r (intrinsic rank).

$$\therefore r(d+k) \ll d \times k.$$

Example

$$\text{if } d = 10000 \\ k = 1000$$

$$\text{using LORA with rank} = 2 \\ 2(10000 + 1000) = 40000$$

which is less than original 10^6 parameters.

LOW RANK DECOMPOSITION

say W_0 is a $d \times k$ matrix

i.e. has $d \times k$ trainable parameters

ΔW s introduced by LORA where

B is a $d \times r$ matrix

A is a $r \times k$ matrix.

*The product BA constructs an ΔW matrix which has same shape as W_0 but with much fewer trainable parameters.