

hw 6

Chen Yuanteng

3039725444

1. Debugging DNNs.

(a)

potential reasons:

c1) Overfitting: 56-layer ^{model} is more complex and deeper, might be more prone to overfitting the training data, while the smaller 20-layer model may have better generalization capabilities.

c2) Vanishing or exploding gradients:

The deeper network architecture could lead to the problem of vanishing or exploding gradients.

c3). maybe training data is limited, it needs more data to train a model as deep and large as 56-layer model.

To mitigate this problem:

c1). Add regularization and dropout layer to reduce the risk of overfitting.

(2) Add residual connections. which can aid information to flow in deeper network and alleviate the issue of vanishing or exploding gradient.

Cb2

the model with layer normalization will not pass the test.

Because layer normalization computes the mean and var across the spatial dimensions for each individual sample in the batch.

In the provided gradient accumulation algorithm, the model accumulates gradients and updates the parameters every accumulated steps.

However, since layer normalization compute statistics per sample, the accumulated gradients from different samples within the same effective batch would have different mean and var values.

(C2). in for loop, `optimizer.zero_grad()` should be implemented after `optimizer.step()`. Otherwise, gradients would be accumulated

during every (inputs, label).

2. Tensor Rematerialization.

(a) to compute activation of layer-9,

we need to compute activations of 6, 7, 8 first, so it needs 4 fwd in total

So to compute activations of 6, 7, 8, 9.

it needs $4 + 3 + 2 + 1 = 10$, same as layer 1-4

$\therefore 2 \times 10 = 20$ fwd in total.

(b). when computing activations of 6, 7, 8, 9.

4 loadmem are necessary,

to compute 1, 2, 3, 4. another 4 loadmem

are needed. So $2 \times 4 = 8$ loadmem in total.

(c)

during a single backward:

in tensor rematerialization,

$$20 \cdot 20 \text{ ns} + 8 \cdot 10 \text{ ns} = 480 \text{ ns.}$$

in storing all activations on this disk:

$$10 \cdot (\text{time of load disk}) = 480 \text{ ns.}$$

$$\text{time of load disk} = 48 \text{ ns}$$

3. Graph Dynamics

$$(a) - G_{t+1} = A \cdot G_t \quad G_0 = E$$

$\therefore G_k = A^k$ and j -th node in G_k is the j -th row of G_k .

\therefore the output of the j -th node at layer k in this network $\rightarrow A_j^k$

(b). using induction:

$$L_0(i,j) = |i=j|$$

$L_1(i,j)$ is the definition of matrix A

assume $L_h(i,j) = [A_h]_{i,j}$ ($h \geq 1$)

try to verify $L_{h+1}(i,j) = [A_{h+1}]_{i,j}$

from node i to node j with distance $= h+1$

\Leftrightarrow from node i to node x with distance $= h$
+ from node x to node j with an edge.

$$\therefore L_{h+1}(i,j) = \sum_{x \in V(G)} L_h(i,x)$$

$$\therefore L_{h+1}(i,j) = \sum_{x=1}^n L_h(i,x) \cdot A_{x,j}$$

$$\therefore L_h(i,j) = [A_h]_{i,j}$$

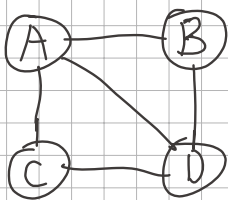
$$\therefore L_{h+1}(i,j) = \sum_{x=1}^n [A_h]_{i,x} A_{x,j} = [A_{h+1}]_{i,j}$$

cc2. in the matrix, each row \Leftrightarrow each node
the graph is a set of nodes:

$$V_j = \sum_{i \in V(j)} V_i$$

cd> replace sum aggregation with max
aggregation:

$$\text{outputs of node } j \text{ at layer } k = \begin{cases} 1 & \text{there is a path from } i \text{ to } j \text{ with length } k \\ 0 & \text{otherwise} \end{cases}$$



	A	B	C	D
A	0	1	1	1
B	1	0	0	1
C	1	0	0	1
D	1	1	1	0

$$A^2 = \begin{matrix} & \begin{matrix} A & B & C & D \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{bmatrix} 3 & 1 & 1 & 2 \\ 1 & 2 & 2 & 1 \\ 1 & 2 & 2 & 1 \\ 2 & 1 & 1 & 2 \end{bmatrix} \end{matrix}$$

