

DFTnet: efficiently training large neural networks

Liu Jiacheng

May 29, 2018

1 Alternative Dense Layer

In a canonical dense layer in neural network, suppose its input size is p and output size is q , then it is described by matrix

$$\begin{bmatrix} w_{11} & w_{12} & \dots & w_{1p} \\ w_{21} & w_{22} & \dots & w_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ w_{q1} & w_{q2} & \dots & w_{qp} \end{bmatrix}$$

the forward propagation is

$$Z = WA$$

and the back propagation is

$$\begin{aligned} dW &= \frac{\partial J}{\partial W} = \frac{1}{m} dZ A^T \\ dA &= \frac{\partial J}{\partial A} = W^T dZ \end{aligned}$$

Now assume $p = q = n$ (we will consider $p \neq q$ later). If the weight matrix W has form

$$\begin{bmatrix} w_1 & w_2 & \dots & w_n \\ w_2 & w_3 & \dots & w_1 \\ \vdots & \vdots & \ddots & \vdots \\ w_n & w_1 & \dots & w_{n-1} \end{bmatrix}$$

Then we would conveniently have (see derivation below) forward propagation

$$Z = w * A = \mathcal{F}^{-1}(\mathcal{F}(w) \cdot \mathcal{F}(A))$$

and back propagation

$$\begin{aligned} dw &= \mathcal{F}\left(\frac{1}{m} \mathcal{F}^{-1}(dZ) \cdot \mathcal{F}(A)\right) \\ dA &= \mathcal{F}(\mathcal{F}^{-1}(dZ) \cdot \mathcal{F}(w)) \end{aligned}$$

where $*$ is linear convolution and $w = [w_1 \ w_2 \ \dots \ w_n]^T$.

2 I/O Size Mismatch

If $p \neq q$, let $n = \max(p, q)$, the weight matrix should be

$$\begin{bmatrix} w_1 & w_2 & \dots & w_p \\ w_2 & w_3 & \dots & w_{(p+1)\%n} \\ \vdots & \vdots & \ddots & \vdots \\ w_q & w_{(q+1)\%n} & \dots & w_{(p+q-1)\%n} \end{bmatrix}$$

In this case, we can still perform DFT by making the following changes:

- Before DFT, pad the input vector with zeros to the end, up to size n .
- After DFT, truncate the output vector from the end, down to size q .

It can be proved that this is equivalent to simple convolution with mismatched I/O size.

3 Time Complexity

Let $n = \max(p, q)$. For each dense layer, the cost of forward propagation is $\Theta(3n \log n + n)$ and the cost of back propagation is $\Theta(3n \log n + 2n)$. This is significantly lower than the $\Theta(pq)$ complexity in canonical dense layer. This implies that we can build wider dense layers given that we reuse the weights circularly.

4 Derivation

$$\begin{aligned} \mathcal{F}(A) &= DFT(A) \\ \mathcal{F}(w) &= DFT(w) \\ \mathcal{F}(Z) &= \mathcal{F}(w)\mathcal{F}(A) \\ Z &= DFT^{-1}(\mathcal{F}(Z)) \end{aligned}$$

$$\begin{aligned} d\mathcal{F}(Z) &= \frac{\partial J}{\partial \mathcal{F}(Z)} = DFT^{-1}(dZ) = \mathcal{F}^{-1}(dZ) \\ d\mathcal{F}(w) &= \frac{\partial J}{\partial \mathcal{F}(w)} = \frac{1}{m} d\mathcal{F}(Z)\mathcal{F}(A) \\ d\mathcal{F}(A) &= \frac{\partial J}{\partial \mathcal{F}(A)} = d\mathcal{F}(Z)\mathcal{F}(w) \\ dw &= \frac{\partial J}{\partial w} = DFT(d\mathcal{F}(w)) = \mathcal{F}(d\mathcal{F}(w)) = \mathcal{F}\left(\frac{1}{m}\mathcal{F}^{-1}(dZ)\mathcal{F}(A)\right) \\ dA &= \frac{\partial J}{\partial A} = DFT(d\mathcal{F}(A)) = \mathcal{F}(d\mathcal{F}(A)) = \mathcal{F}(\mathcal{F}^{-1}(dZ)\mathcal{F}(w)) \end{aligned}$$

5 Experiment

I trained a DFTnet: a neural network composed purely by DFT dense layers and ReLU activation. It has layer width 785, 4096, 1024, 256, 64, 10. The model is fit into MNIST dataset and the result is comparable to a canonical neural network.

