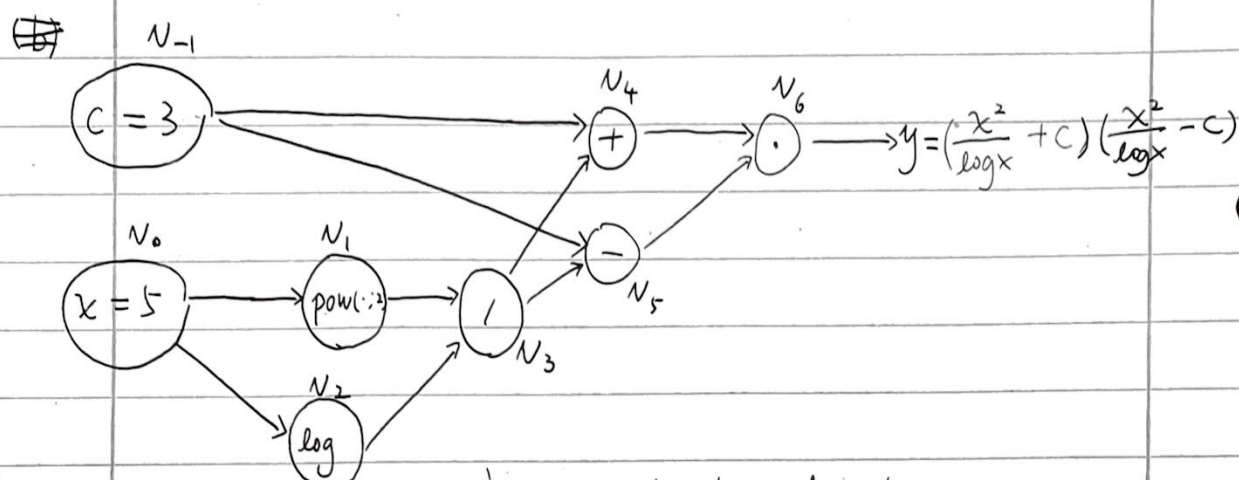


$$1. (a) \frac{\partial \left[\left(\frac{x^2}{\log x} + c \right) \left(\frac{x^2}{\log x} - c \right) \right]}{\partial x}$$

$$= \frac{\partial \left[\left(\frac{x^2}{\log x} + c \right) \left(\frac{x^2}{\log x} - c \right) \right]}{\partial \left(\frac{x^2}{\log x} + c \right)} \cdot \frac{\partial \left(\frac{x^2}{\log x} + c \right)}{\partial \left(\frac{x^2}{\log x} \right)} \left[\frac{\partial \left(\frac{x^2}{\log x} \right)}{\partial (x^2)} \cdot \frac{\partial (x^2)}{\partial x} + \frac{\partial \left(\frac{x^2}{\log x} \right)}{\partial (\log x)} \cdot \frac{\partial (\log x)}{\partial x} \right]$$

$$+ \frac{\partial \left[\left(\frac{x^2}{\log x} + c \right) \left(\frac{x^2}{\log x} - c \right) \right]}{\partial \left(\frac{x^2}{\log x} - c \right)} \cdot \frac{\partial \left(\frac{x^2}{\log x} - c \right)}{\partial \left(\frac{x^2}{\log x} \right)} \left[\frac{\partial \left(\frac{x^2}{\log x} \right)}{\partial (x^2)} \cdot \frac{\partial (x^2)}{\partial x} + \frac{\partial \left(\frac{x^2}{\log x} \right)}{\partial (\log x)} \cdot \frac{\partial (\log x)}{\partial x} \right]$$



(b) Forward trace:

$$\begin{aligned} N_{-1} = c &= 3 \\ N_0 = x &= 5 \\ N_1 = N_0^2 &= 25 \\ N_2 = \log N_0 &= 0.699 \\ N_3 = N_1 / N_2 &= 35.8 \\ N_4 = N_{-1} + N_3 &= 38.8 \\ N_5 = N_3 - N_{-1} &= 32.8 \\ N_6 = N_4 \cdot N_5 &= 1272.6 \\ y = N_6 &= 1272.6 \end{aligned}$$

(c) backward trace:

$$\begin{aligned} \bar{v}_6 &= \bar{y} = 1 \\ \bar{v}_5 &= \bar{v}_6 \frac{\partial v_6}{\partial v_5} = \bar{v}_6 \times v_4 = 1 \times v_4 = 38.8 \\ \bar{v}_4 &= \bar{v}_6 \frac{\partial v_6}{\partial v_4} = \bar{v}_6 \times v_5 = 1 \times 32.8 = 32.8 \\ \bar{v}_3 &= \bar{v}_4 \frac{\partial v_4}{\partial v_3} + \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_4 \times 1 + \bar{v}_5 \times (-1) = 71.6 \\ \bar{v}_2 &= \bar{v}_3 \frac{\partial v_3}{\partial v_2} = \bar{v}_3 \times \left(-\frac{v_1}{v_2^2} \right) = -3663.5 \\ \bar{v}_1 &= \bar{v}_3 \frac{\partial v_3}{\partial v_1} = \bar{v}_3 \times \frac{1}{v_2} = 102.4 \\ \bar{v}_0 &= \bar{v}_1 \frac{\partial v_1}{\partial v_0} + \bar{v}_2 \frac{\partial v_2}{\partial v_0} = \bar{v}_1 \cdot (2v_0) + \bar{v}_2 \cdot \frac{1}{v_0} = 291.3 \\ \bar{v}_{-1} &= \bar{v}_4 \frac{\partial v_4}{\partial v_{-1}} + \bar{v}_5 \frac{\partial v_5}{\partial v_{-1}} = \bar{v}_4 \times 1 + \bar{v}_5 \times (-1) = -6 \\ \bar{x} &= \bar{v}_0 = 291.3 \\ \bar{c} &= \bar{v}_{-1} = -6 \end{aligned}$$

We only need to do some simple calculation in every step of backprop, which is simpler for computers.

2. (a) On iteration t :

$$m^t = \beta m^{t-1} + (1-\beta)g^t, \quad \text{calculate the accumulated momentum}$$

$$v^t = \gamma v^{t-1} + (1-\gamma)(g^t)^2, \quad \dots \dots \dots \text{second moment}$$

$$\hat{m}^t = \frac{m^t}{1-(\beta)^t}, \quad \hat{v}^t = \frac{v^t}{1-(\gamma)^t} \quad \text{correction of initial bias}$$

$$w^{t+1} = w^t - \alpha \frac{\hat{m}^t}{\sqrt{\hat{v}^t + \epsilon}} \quad \text{update weight}$$

$$m^0 = 0, \quad v^0 = 0.$$

usually $\beta \sim 0.99$, $\gamma \sim 0.999$, $\epsilon \sim 10^{-8}$

$$(b) \quad m^1 = 0 \cdot \beta + (1-\beta)g^1 = (1-\beta)g^1, \quad \hat{m}^1 = \frac{m^1}{1-\beta} = g^1.$$

$$v^1 = \gamma \cdot 0 + (1-\gamma)(g^1)^2 = (1-\gamma)(g^1)^2, \quad \hat{v}^1 = \frac{v^1}{1-\gamma} = (g^1)^2.$$

$$w^2 = w^1 - \alpha \frac{g^1}{|g^1| + \epsilon} \approx w^1 - \alpha \text{sign}(g^1).$$

$$(c) \quad m^2 = (1-\beta)\beta g^1 + (1-\beta)g^2, \quad \hat{m}^2 = \frac{m^2}{1-\beta^2} = (\beta g^1 + g^2)/(1+\beta)$$

$$v^2 = (1-\gamma)\gamma(g^1)^2 + (1-\gamma)(g^2)^2, \quad \hat{v}_2 = \frac{v^2}{1-\gamma^2} = [\gamma(g^1)^2 + (g^2)^2]/(1+\gamma)$$

$$\alpha \frac{\hat{m}^2}{\sqrt{\hat{v}^2 + \epsilon}} \approx \alpha \frac{\beta g^1 + g^2}{\sqrt{\gamma(g^1)^2 + (g^2)^2}} \frac{\sqrt{1+\gamma}}{1+\beta} \quad \left(\begin{array}{l} \hat{m} \text{ is the weighted average of gradient} \\ \hat{v} \text{ is the weighted squared average.} \end{array} \right)$$

(d) ~~Initialize $m^0 =$~~

(e) L2-Regularization in Loss function:

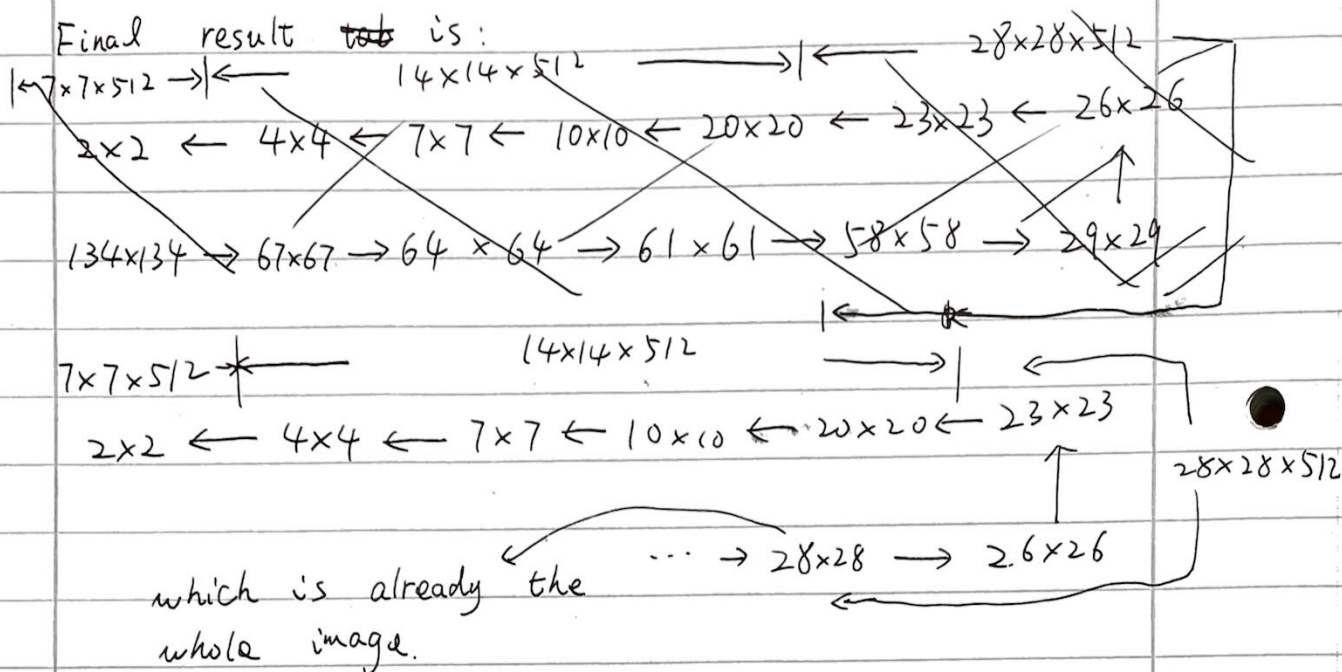
$$L \rightarrow L' = L + \lambda \|w\|_2$$

$$\Rightarrow g^t \rightarrow g'^t = g^t + 2\lambda w^t,$$

this

3.(a) A perceptron in $7 \times 7 \times 512$ maxpooling layer depends on ~~2x2~~ previous output, which depends on a $(3+1) \times (3+1)$ region of previous layer, which again depends on $(3+3+1) \times (3+3+1)$ region of previous layer, and soon.

Final result ~~is~~ is:



so the FOV is 224×224 .

(b) # parameters = ~~3x3x64~~

(with bias)
in Convolutional
layers

$$\begin{aligned}
 & (3 \times 3 \times 3 + 1) \times 64 + (3 \times 3 \times 64 + 1) \times 64 \\
 & + (3 \times 3 \times 64 + 1) \times 128 + (3 \times 3 \times 128 + 1) \times 128 \\
 & + (3 \times 3 \times 128 + 1) \times 256 + (3 \times 3 \times 256 + 1) \times 256 \times 2 \\
 & + (3 \times 3 \times 256 + 1) \times 512 + (3 \times 3 \times 512 + 1) \times 512 \times 2 \\
 & + \cancel{(3 \times 3 \times 512 + 1) \times 512 \times 3} + (3 \times 3 \times 512 + 1) \times 512 \times 3
 \end{aligned}$$

parameters in
fully connected
layers

$$\begin{aligned}
 & = (7 \times 7 \times 512 \times 4096 + 4096) \\
 & + 4096 \times 4096 + 4096 \\
 & + 4096 \times 1000 + 1000 \\
 & = 123642856 \approx 1.24 \times 10^8
 \end{aligned}$$

parameters in total = $138357544 \approx 1.38 \times 10^8$

$$\frac{\text{\# of conv} \times \text{\# of convolutional}}{\text{\# of fully connected}} = 11.90\%$$