

# backprop

December 8, 2024

## 1. Reverse Mode Automatic Differentiation

[ ]: d)

```
[11]: import torch

x = torch.tensor(3.0, requires_grad=True)
c = torch.tensor(5.0, requires_grad=True)

log_x = torch.log(x)
div = (x**2) / log_x
mult = (div + c) * (div - c)

mult.backward()

grad_x = x.grad
grad_c = c.grad

print(f"Output: {mult.item()}")
print(f"Gradient w.r.t x: {grad_x.item()}")
print(f"Gradient w.r.t c: {grad_c.item()}")
```

Output: 42.11137008666992

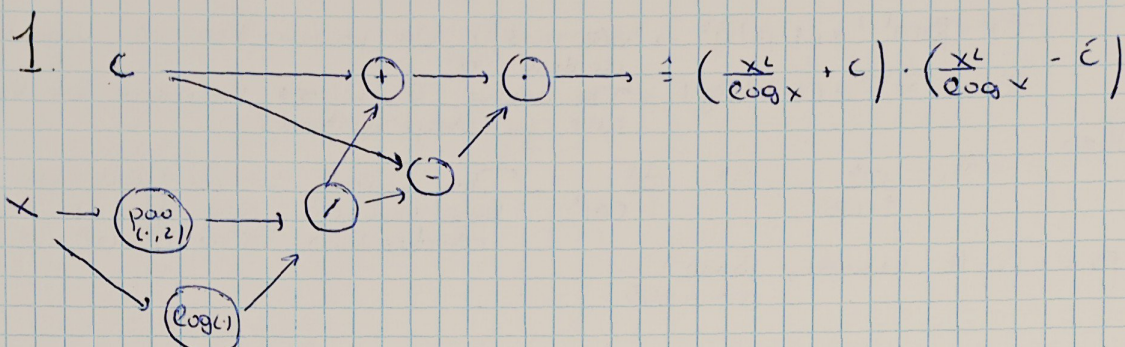
Gradient w.r.t x: 48.756893157958984

Gradient w.r.t c: -10.0

[ ]:

# Exercise Sheet 7

Arsen Muramatov  
Pylyp Filippov

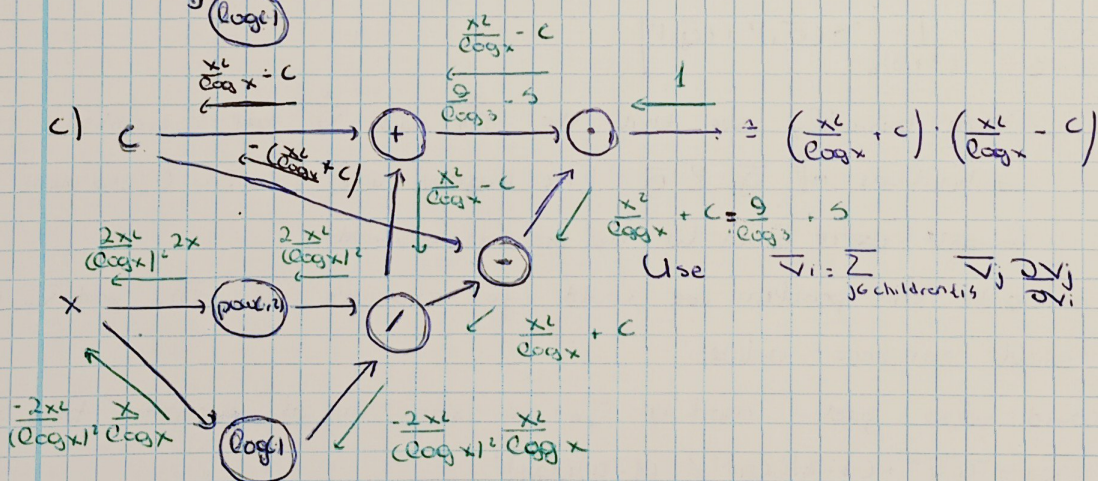
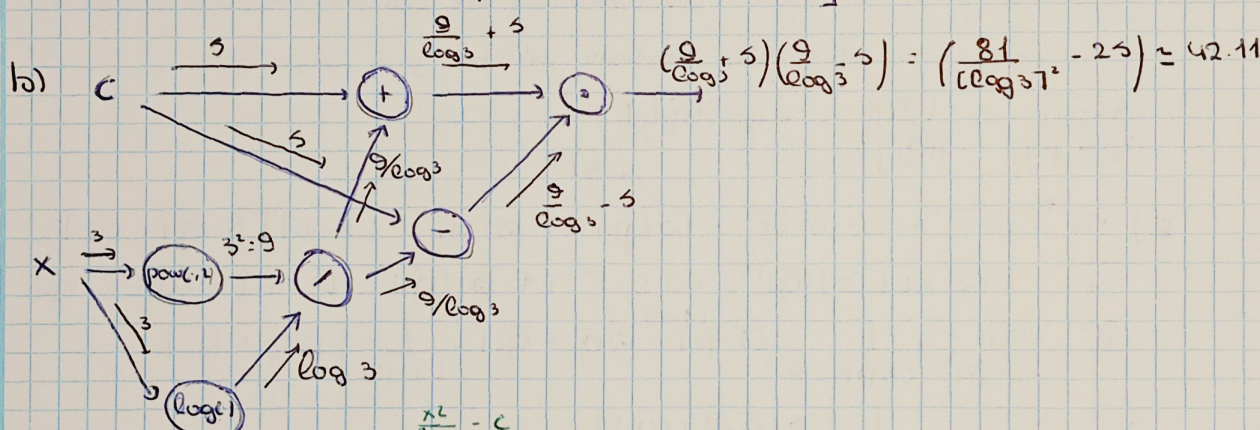


a) 
$$d \left[ \left( \frac{x^c}{\log x} + c \right) \left( \frac{x^c}{\log x} - c \right) \right] =$$

$$= \left[ \frac{\partial \left[ \frac{x^c}{\log x} \right]}{\partial x^c} \frac{dx^c}{dx} + \frac{\partial \left[ \frac{x^c}{\log x} \right]}{\partial \log x} \frac{d \log x}{dx} \right] x$$

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

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



$$\frac{d}{dx} \left[ \left( \frac{x^c}{\log x} + c \right) \left( \frac{x^c}{\log x} - c \right) \right] = \frac{2x^c}{(\log x)^2} [2x \log x - x] \Big|_{x=3} = \frac{18}{(\log 3)^2} [6 \log 3 - 3] \approx 48.76$$

$$\frac{d}{dc} \left[ \left( \frac{x^c}{\log x} - c \right) \left( \frac{x^c}{\log x} + c \right) \right] = \frac{x^c}{\log x} - c - \frac{x^c}{\log x} - c = -2c \Big|_{c=5} = -10$$

~~I would say that in this case symbolic differentiation is~~

If we compare b) with the result in a) we see that c) might look a bit easier. But at the end, we're basically doing the same thing (chain rule).



2. a)  $\bar{w}^{t+1} = \bar{w}^t - \frac{\alpha \hat{m}^t}{\sqrt{\hat{v}^t} + \epsilon} \rightarrow$  In each iteration, we move the parameters in certain direction multiplied by the learning rate

$\bar{m}^t = \beta \bar{m}^{t-1} + (1-\beta) \bar{g}^t \rightarrow$  We modify the previous direction using the gradient  $\bar{g}^t$

$\hat{v}^t = \delta \hat{v}^{t-1} + (1-\delta) (\bar{g}^t)^2 \rightarrow$  The vector  $\hat{v}^t$  is used to reduce the learning rate for steep gradients

Also

$\hat{m}^t = \frac{\bar{m}^t}{1 - (\beta)^t}, \hat{v}^t = \frac{\hat{v}^t}{1 - (\delta)^t} \rightarrow$  This is used to counteract the bias towards zero for the initialization  $\bar{m}^0 = 0, \hat{v}^0 = 0$

$\beta, \delta \in [0, 1]$ , Usually  $\beta \sim 0.9$   
 $\delta \sim 0.999$

b) First iteration

$$\left. \begin{aligned} \bar{m}^1 &= \beta \bar{m}^0 + (1-\beta) \bar{g}^1 = (1-\beta) \bar{g}^1 \\ \hat{v}^1 &= \delta \hat{v}^0 + (1-\delta) (\bar{g}^1)^2 = (1-\delta) (\bar{g}^1)^2 \\ \hat{m}^1 &= \frac{\bar{m}^1}{1-\beta} = \frac{1}{1-\beta} (1-\beta) \bar{g}^1 = \bar{g}^1 \\ \hat{v}^1 &= \frac{(1-\delta) (\bar{g}^1)^2}{1-\delta} = (\bar{g}^1)^2 \end{aligned} \right\} \frac{\hat{m}^1}{\sqrt{\hat{v}^1} + \epsilon} = \frac{\bar{g}^1}{\sqrt{(\bar{g}^1)^2} + \epsilon} = \frac{\bar{g}^1}{|\bar{g}^1| + \epsilon}$$

$$\frac{\hat{m}^1}{\sqrt{\hat{v}^1} + \epsilon} = \frac{\bar{g}^1}{|\bar{g}^1| + \epsilon} \approx \frac{\bar{g}^1}{|\bar{g}^1|} = \frac{\bar{g}^1}{\sqrt{(\bar{g}^1)^2}}$$

Bearing in mind that the operations are element-wise, we conclude that the components of the gradient will be just  $\text{sign}(\bar{g}_i)$

$$\frac{(\bar{g}_i)_j}{|\bar{g}_i|} = \frac{(\bar{g}_i)_j}{|(\bar{g}_i)_j|} = \text{sign}[(\bar{g}_i)_j]$$

$$c) \bar{m}^2 = \beta \bar{m}^1 + (1-\beta) \bar{g}^2 = \beta \bar{g}^1 + (1-\beta) \bar{g}^2 = [1-\beta] [\beta \bar{g}^1 + \bar{g}^2]$$

$$\hat{v}^2 = \delta \hat{v}^1 + (1-\delta) (\bar{g}^2)^2 = \delta (\bar{g}^1)^2 + (1-\delta) (\bar{g}^2)^2 = [1-\delta] [\delta (\bar{g}^1)^2 + (\bar{g}^2)^2]$$

$$\hat{m}^2 = \frac{\bar{m}^2}{1-\beta^2} = \frac{1}{(1-\beta)(1-\beta)} [1-\beta] [\beta \bar{g}^1 + \bar{g}^2] = \frac{1}{(1-\beta)} [\beta \bar{g}^1 + \bar{g}^2]$$

$$\hat{v}^2 = \frac{\hat{v}^2}{1-\delta^2} = \frac{1}{(1-\delta)} [\delta (\bar{g}^1)^2 + (\bar{g}^2)^2]$$

$\frac{\hat{m}^2}{\sqrt{\hat{v}^2} + \epsilon} \rightarrow$  We basically see that the  $\hat{m}^2$  and  $\hat{v}^2$  are weighted combinations of  $\bar{g}_1$  &  $\bar{g}_2$ . The bigger are  $\beta$  and  $\delta$ , the more weight we're given to the previous gradient

d) To mitigate this problem we could for example set  $\bar{m}_0$  and  $\hat{v}_0$  to some small non-zero values

We could also just modify  $\hat{v}_0$  to some initial constant big enough

$$\text{so } \hat{v}_1 = \delta \hat{v}^0 + (1-\delta) (\bar{g}_1)^2 \neq (1-\delta) (\bar{g}_1)^2$$

$$\text{basically } |\delta \hat{v}^0| \sim |(1-\delta) (\bar{g}_1)^2|$$

e) Yes, there is a difference

$$\left. \begin{aligned} \mathcal{L}_{\text{tot}}(w) &= \mathcal{L}_{\text{orig}}(w) + \lambda \|w\|_2^2 \\ \frac{\partial \mathcal{L}_{\text{tot}}}{\partial w} &= \frac{\partial \mathcal{L}_{\text{orig}}}{\partial w} + 2\lambda w \end{aligned} \right\} \text{The penalty explicitly contributes to the gradient.}$$

$\bar{w}^{t+1} = \bar{w}^t - \frac{\alpha \hat{m}^t}{\sqrt{\hat{v}^t} + \epsilon} - \alpha \lambda \bar{w}^t \rightarrow$  In this case we scale down the weight during each step, independently of the gradient.

We can conclude that the second option is better (weight decay applied to weights) because it is applied uniformly across all parameters, while ~~the~~ in the  $\mathcal{L}_2$  regularization the reg. term is affected by Adam's adaptive learning rates. Basically it is better to decouple the weight decay from gradient update



### 3 Receptive Field of VGG16

(a)

$$r_l = (r_{l+1} - 1) * s_{l+1} + k_{l+1}$$

Where

- $r_l$  – the size of the receptive field of a pixel on layer  $l$
- $s_l$  – stride
- $k_l$  – kernel size

Convolutional layers in VGG16 have  $k = 3$ ,  $s = 1$ , while max pooling layers can be represented as convolutional layers with  $k = 2$ ,  $s = 2$  with constant weights.

Layer	Receptive field
18. MaxPool (Output):	$r_{18} = 1$
17. Conv:	$r_{17} = (r_{18} - 1) * 2 + 2 = r_{18} * 2 = 2$
16. Conv:	$r_{16} = (r_{17} - 1) * 1 + 3 = r_{17} + 2 = 4$
15. Conv:	$r_{15} = 6$
14. MaxPool:	$r_{14} = 8$
13. Conv:	$r_{13} = 16$
12. Conv:	$r_{12} = 18$
11. Conv:	$r_{11} = 20$
10. MaxPool:	$r_{10} = 22$
9. Conv:	$r_9 = 44$
8. Conv:	$r_8 = 46$
7. Conv:	$r_7 = 48$
6. MaxPool:	$r_6 = 50$
5. Conv:	$r_5 = 100$
4. Conv:	$r_4 = 102$
3. MaxPool:	$r_3 = 104$
2. Conv:	$r_2 = 208$
1. Conv:	$r_1 = 210$
0. Input:	$r_0 = 212$

(b)

```
s = lambda x, y: x**2 + y
n_params_conv = 3 * (64 * 2 + 128 * 2 + 256 * 3 + 512 * 3 + 512 * 3)
n_params_fc = (512 * 7 * 7) * (4096) + 4096**2 + 4096 * 1000
print("Total parameters:", "%.5E" % (n_params_conv + n_params_fc))
print("Ratio:", "%.5E" % (n_params_conv / n_params_fc))
```

Total parameters: 1.23646E+08 Ratio: 1.02496E-04