Names: Philipp Köhler, Alexander Bespalov

**1**

**a)**

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

$$a(x) = \left(\frac{x^2}{\log(x)} + c\right)$$

$$b(x) = \left(\frac{x^2}{\log(x)} - c\right)$$

$$c(x) = \frac{x^2}{\log(x)}$$

$$\frac{\partial f(x)}{\partial x} = \frac{\partial f(x)}{\partial a(x)} \frac{\partial a(x)}{\partial c(x)} \left(\frac{\partial c(x)}{\partial x^2} \frac{\partial x^2}{\partial x} + \frac{\partial c(x)}{\partial \frac{1}{\ln(x)}} \frac{\partial \frac{1}{\ln(x)}}{\partial \ln(x)} \frac{\partial \ln(x)}{x}\right)$$

$$+ \frac{\partial f(x)}{\partial b(x)} \frac{\partial b(x)}{\partial c(x)} \left(\frac{\partial c(x)}{\partial x^2} \frac{\partial x^2}{\partial x} + \frac{\partial c(x)}{\partial \frac{1}{\ln(x)}} \frac{\partial \frac{1}{\ln(x)}}{\partial \ln(x)} \frac{\partial \ln(x)}{x}\right)$$

$$= \left(\frac{\partial f(x)}{\partial a(x)} \frac{\partial a(x)}{\partial c(x)} + \frac{\partial f(x)}{\partial b(x)} \frac{\partial b(x)}{\partial c(x)}\right)\left(\frac{\partial c(x)}{\partial x^2} \frac{\partial x^2}{\partial x} + \frac{\partial c(x)}{\partial \frac{1}{\ln(x)}} \frac{\partial \frac{1}{\ln(x)}}{\partial \ln(x)} \frac{\partial \ln(x)}{\partial x}\right)$$

**b)**

$$x = 3, \quad c = 5$$

$$x^2 = 9, \quad \ln(x) = \ln(3) \approx 1{,}1$$

$$c(3) = \frac{9}{\ln(3)} \approx 8{,}2$$

$$a(3) = \frac{9}{\ln(3)} + 5 \approx 13{,}2, \quad b(3) = \frac{9}{\ln(3)} - 5 \approx 3{,}2$$

$$f(3) \approx 42{,}1$$

**c)**

$$\left.\frac{\partial \ln(x)}{\partial x}\right|_{x=3} = \left.\frac{1}{x}\right|_{x=3} = \frac{1}{3}$$

$$\left.\frac{\partial \frac{1}{\ln(x)}}{\partial \ln(x)}\right|_{\ln(x) = \ln(3) \approx 1{,}1} = -\frac{1}{\ln(3)^2} \approx -0{,}8$$

$$\left.\frac{\partial c(x)}{\partial \frac{1}{\ln(x)}}\right|_{x^2 = 9} = \left.x^2\right|_{x^2 = 9} = 9$$

<span style="color:red">Aren't you just evaluating the chain rule here?</span>

$$\left.\frac{\partial x^2}{\partial x}\right|_{x=3} = \left.2x\right|_{x=3} = 6$$

$$\left.\frac{\partial c(x)}{\partial x^2}\right|_{\frac{1}{\ln(x)} = \frac{1}{\ln(3)} \approx 0{,}9} = \left.\frac{1}{\ln(x)}\right|_{\frac{1}{\ln(x)} = \frac{1}{\ln(3)} \approx 0{,}9} = 0{,}9$$

$$\frac{\partial b(x)}{\partial c(x)} = 1, \quad \frac{\partial a(x)}{\partial c(x)} = 1$$

$$\left.\frac{\partial f(x)}{\partial b(x)}\right|_{a(x) = a(3) \approx 13{,}2} = \left.a(x)\right|_{a(x) = a(3) \approx 13{,}2} \approx 13{,}2$$

$$\left.\frac{\partial f(x)}{\partial a(x)}\right|_{s(x)=b(3)=3,2} = b(x)\Big|_{b(x)=b(3)=3,2} \approx 3,2$$

$$\Rightarrow \left.\frac{\partial f(x)}{\partial x}\right|_{x=3,\, c=5} \approx \left(3,2 \cdot 1 + 13,2 \cdot 1\right) \cdot \left(0,5 \cdot 6 + 9 \cdot (-0,8) \cdot \frac{1}{3}\right) = 45,2$$

In this case and by hand symbolic differentiation is easier, because a lot of annotation disappears.

d)

## 1 Reverse Mode Automatic Differentiation

d)

```python
import torch

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

f1 = x**2
f2 = torch.log(x)
f3 = f1 / f2
f4 = f3 + c
f5 = f3 - c
output = f4 * f5

output.backward()

dx = x.grad
dc = c.grad

print(f"Derivative with respect to x: {dx}")
print(f"Derivative with respect to c: {dc}")
```

```
Derivative with respect to x: 48.756893157958984
Derivative with respect to c: -10.0
```

solution is different from calculated x with 0.5 distance, which is due to the rounding of intermediate steps in the calculation by hand.

**2**

**a)**

$$w^{t+1} = w^t - \alpha \frac{\hat{m}^t}{\sqrt{\hat{v}^t} + \varepsilon}$$

Update of the parameters with learningrate $\alpha$.

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

$m^t$ is the momentum that smooths the gradient $g^t$ with the hyperparameter $\beta \in [0,1]$

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

$v^t$ reduces the gradient for steep gradients

The division of $(1-(\beta)^t)$ and $(1-(\gamma)^t)$ counteracts the initialisation bias towards 0.

**b)**

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

$$\hat{m}^1 = (\beta m^0 + (1-\beta) g^1)/(1-(\beta)^1) = g$$

$$\hat{v}^1 = (\gamma v^0 + (1-\gamma)(g^1)^2)/(1-(\gamma)^1) = g^2$$

$$\frac{\hat{m}^1}{\sqrt{\hat{v}^1} + \varepsilon} = \frac{g}{\sqrt{g^2} + \varepsilon} = \frac{g}{|g|} = \text{sign}(g)$$

**c)**

$$m^1 = (1-\beta) g^1 \quad , \quad v^1 = (1-\gamma)(g^1)^2$$

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

$$\hat{v}^2 = (\gamma v^1 + (1-\gamma)(g^2)^2)/(1-(\gamma)^2) = \frac{\gamma(1-\gamma)}{(1-\gamma^2)}(g^1)^2 + \frac{1-\gamma}{(1-\gamma^2)}(g^2)^2$$

$$\frac{\hat{m}^2}{\sqrt{\hat{v}^2} + \varepsilon} = \frac{\frac{\beta(1-\beta)}{(1-\beta^2)} g^1 + \frac{(1-\beta)}{(1-\beta^2)} g^2}{\sqrt{\frac{\gamma(1-\gamma)}{(1-\gamma^2)}(g^1)^2 + \frac{1-\gamma}{(1-\gamma^2)}(g^2)^2} + \varepsilon}$$

**d)**

Smaller initial learningrates in the initial steps (learningrate warmup) can solve this issue of the dominating sign(g).

**e)**

Adam has an adaptive learningrate therefore L2 regularization is not the same as weight decay. In the case of SGD it would be. In Adam the regularization is varied by the adaptive learningrate wich leads to inconsistency and less predictable behaviour. The weight decay is preferred.

**3**

**a)** maxpooling Kernel $k=2$, with stride $s=2$, padding $p=0$

convolution Kernel $k=3$, with stride $s=1$, padding $p=1$

receptive field $r$

$$r_{out} = r_{in} + (k-1)\cdot jump \quad , \quad jump_{out} = jump_{in}\cdot s \quad , \quad \text{general formula:}$$

initial: $r_0 = 1$ , $jump_0 = 1$ $\qquad r_L = 1 + \sum_{i=1}^{L}\prod_{j=1}^{i-1} s_j\cdot(k_i-1)$

Conv 1: $\quad r_1 = r_0 + (k-1)\cdot jump_0 \quad , \quad jump_1 = jump_0 \cdot s$

$\qquad\qquad = 1 + (3-1)\cdot 1 \qquad\qquad = 1\cdot 1$

$\qquad\qquad = 3 \qquad\qquad\qquad\qquad = 1$

Conv 2: $\quad r_2 = 3 + (3-1)\cdot 1 = 5 \quad , \quad jump_2 = 1\cdot 1 = 1$

maxp. 1: $\quad r_3 = 5 + (2-1)\cdot 1 = 6 \quad , \quad jump_3 = 1\cdot 2 = 2$

Conv. 3: $\quad r_4 = 6 + (3-1)\cdot 2 = 10 \quad , \quad jump_4 = 2\cdot 1 = 2$

Conv. 4: $\quad r_5 = 10 + (3-1)\cdot 2 = 14 \quad , \quad jump_5 = 2\cdot 1 = 2$

maxp. 2: $\quad r_6 = 14 + (2-1)\cdot 2 = 16 \quad , \quad jump_6 = 2\cdot 2 = 4$

$$\vdots$$

<span style="color:red">Perfect</span>

$$r_{out} = 212$$

**b)**

Params in conv layer $= \#\text{Filters}\cdot\left[\#\text{input channels}\cdot\text{Kernel size}^{\sigma_{bias}} + 1\right]$

<span style="color:red">^2 as 2-dim. filters are used</span>

maxpool no params

Params in FC layer $= \#\text{inputs}\cdot\#\text{outputs} + \#\text{outputs}$

First layer:

$\quad \#\text{param} = 64\cdot(3\cdot 3 + 1) = 640$

2. layer: $\quad \#\text{param} = 64\cdot(64\cdot 3 + 1) = 12\,352$

$$\vdots$$

$\#$ total params $= 123\,066\,664$

$\#$ fc param $= 123\,642\,856$

$\#$ conv param $= 5\,423\,808$

ratio $= \dfrac{\#\text{conv param}}{\#\text{fc param}} = 0,044$