



$$\frac{\partial f_1(x)}{\partial x} = 2x \quad \frac{\partial f_2(x)}{\partial x} = \frac{7}{x}$$

$$\frac{\partial f_3}{\partial x} = \frac{\partial \left(\frac{f_1(x)}{f_2(x)} \right)}{\partial x} = \frac{\partial f_1(x)}{\partial x} \cdot \frac{7}{f_2(x)} + \frac{7}{f_2^2(x)} \cdot \frac{\partial f_2(x)}{\partial x} \cdot f_1(x)$$

$$\frac{\partial f_4(x)}{\partial x} = \frac{\partial f_3(x)}{\partial x} + \frac{\partial C}{\partial x} \quad \frac{\partial f_5(x)}{\partial x} = \frac{\partial f_3(x)}{\partial x} - \frac{\partial C}{\partial x}$$

$$\frac{\partial f_6(x)}{\partial x} = \frac{\partial f_4(x)}{\partial x} \cdot f_5(x) + \frac{\partial f_5(x)}{\partial x} \cdot f_4(x)$$

$$l) \quad x=3 \quad l=5$$

$$f_1(3) = 9$$

$$f_3(3) = \frac{9}{\log(3)} \approx 8,792$$

$$f_5(3,5) = \left(\frac{9}{\log(3)} \right)^{-5} \approx 3,792$$

$$f_2(3) = \log(3) \approx 1,099$$

$$f_4(3,5) = \left(\frac{9}{\log(3)} \right)^{+5} \approx 73,792$$

$$f_6(3,5) = \left(\frac{9}{\log(3)} \right)^{+5} \left(\frac{9}{\log(3)} \right)^{-5} \approx 42$$

$$1) \quad \frac{df_6}{df_5} = f_4 \approx 73,792$$

$$\frac{df_6}{df_4} = f_5 \approx 3,792$$

$$2) \quad \frac{df_5}{df_3} = 7$$

$$\frac{df_4}{df_3} = 7$$

$$\frac{df_6}{df_3} = \frac{df_6}{df_4} \frac{df_4}{df_3} + \frac{df_6}{df_5} \frac{df_5}{df_3} \approx 76,384$$

$$3) \quad \frac{df_3}{df_1} = \frac{1}{f_2}$$

$$\frac{df_3}{df_2} = \frac{f_1}{f_2^2}$$

$$\frac{df_6}{df_3} \frac{df_3}{df_1} \approx 14,97$$

$$\frac{df_6}{df_3} \frac{df_3}{df_1} \approx -727,63$$

$$\frac{df_1}{dx} = 2x = 6$$

$$\frac{df_2}{dx} = \frac{1}{x} = \frac{1}{3}$$

$$\frac{df_6}{df_1} \frac{df_1}{dx} + \frac{df_6}{df_2} \frac{df_2}{dx} \approx 48,8$$

\Rightarrow backward is easier, no explicit derivations required

$$K2 (1) w^{t+1} = w^t - \alpha \frac{\hat{m}^t}{\sqrt{\hat{v}} + \epsilon}$$

$$(2) m^t = \beta m^{t-1} + (1 - \beta) g^t$$

$$(3) v^t = \gamma v^{t-1} + (1 - \gamma) (g^t)^2$$

$$(4) \hat{m}^t = \frac{m^t}{1 - (\beta)^t}, \quad (5) \hat{v}^t = \frac{v^t}{1 - (\gamma)^t},$$

(1) update of the weights

$\frac{\alpha}{\sqrt{\hat{v}} + \epsilon}$ rescales the learning rate

ϵ avoids division by 0

(2) updates momentum by involving the previous gradient

(3) reduces the learning rate for steep gradients

(4)(5) avoids bias through initialization for $m^0 = 0$ and $v^0 = 0$

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

$$\hat{m}^1 = g^1$$

$$\hat{v}^1 = (g^1)^2$$

$$w^1 = w^0 - \alpha \frac{g^1}{\sqrt{(g^1)^2}} = w^0 - \alpha \left(\frac{g_1^1}{|g_1^1|}, \frac{g_2^1}{|g_2^1|}, \dots \right)^T = w^0 - \alpha \text{sign } g^1$$

↖ component wise

$$c) m^2 = \beta (1 - \beta) g^1 + (1 - \beta) g^2$$

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

$$v^2 = \gamma (1 - \gamma) (g^1)^2 + (1 - \gamma) (g^2)^2$$

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

$$W^2 = W^1 - \gamma \operatorname{sign} g^1 - \gamma \frac{(\beta g^1 + g^2) \sqrt{1+\gamma}}{(1+\beta) \sqrt{\gamma(g^1)^2 + (g^2)^2}}$$

d) Train more than one epoche!

e) Yes, it makes a difference. Including $\|W\|_2^2$ in the loss affects the gradient computation, causing regularization to interact with adam's adaptive learning rates.

In contrast, AdamW applies weight decay directly during updates, decoupling regularization from optimization. This leads to more stable training and better generalization, making AdamW the preferred approach.

Sure

Exercise07_Kusch_Reckermann_Weinreich_code

December 5, 2024

1 7.1

1.1 (d)

```
[2]: import torch
```

```
[4]: 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
    f6 = f4 * f5

    # Backward computation
    f6.backward()

    # Print gradients
    print(f"df6/dx: {x.grad}")
```

```
df6/dx: 48.756893157958984
```