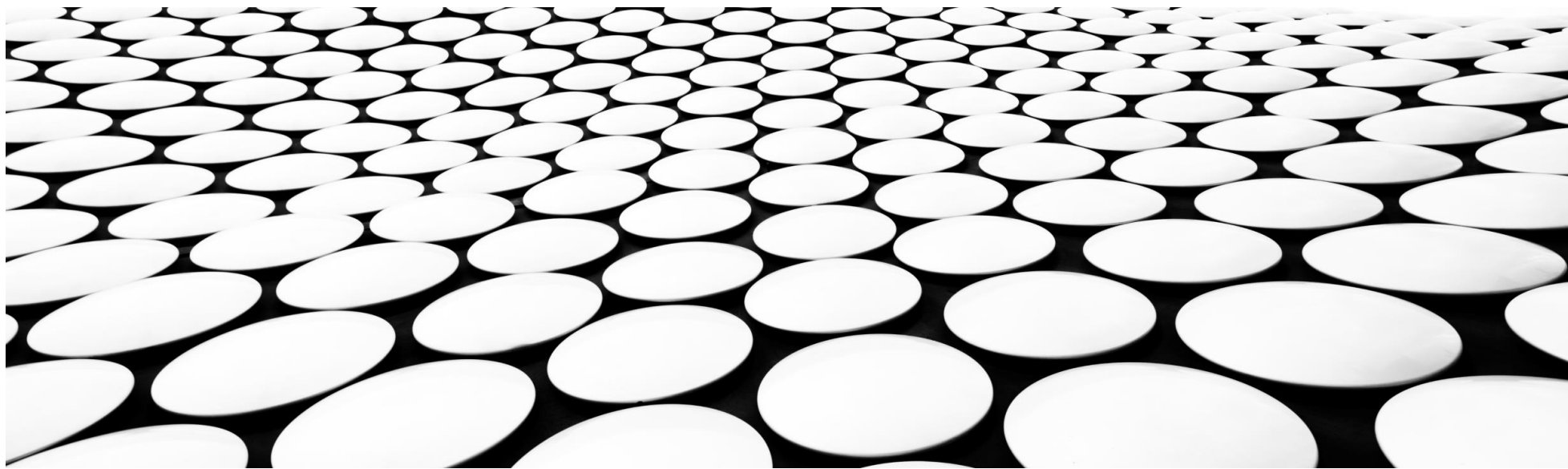


深度学习

邱怡轩



今天的主题

- 反向传播算法（续）
- PyTorch 自动微分

回到神经网络

- 利用之前介绍的方法，考虑一个简单的前馈神经网络

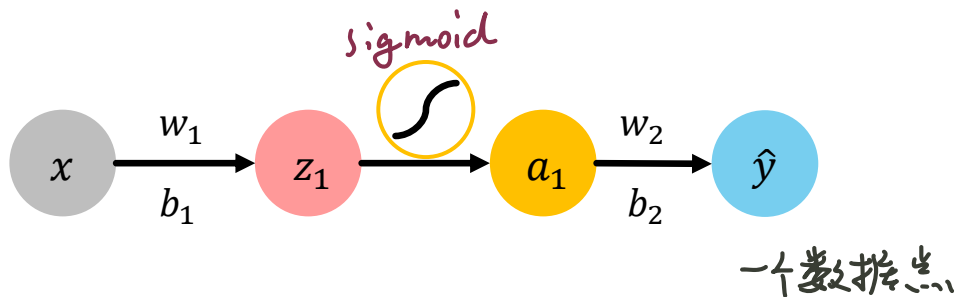
- $a_0 = x, z_1 = w_1 a_0 + b_1, a_1 = \sigma(z_1)$

- $\hat{y} = z_2 = w_2 a_1 + b_2$

- $l = (y - \hat{y})^2$

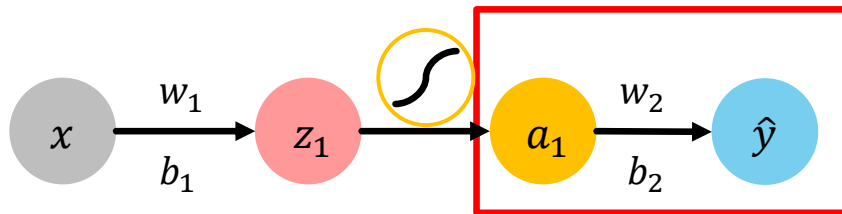
- 需要计算 $\frac{dl}{dw_1}, \frac{dl}{db_1}, \frac{dl}{dw_2}, \frac{dl}{db_2}$

↑
Sigmoid函数



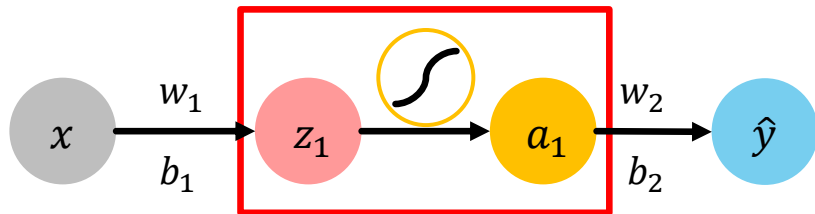
回到 神经网络

- $\hat{y} = z_2$ 的上游导数
- $\frac{dl}{dz_2} = -2(y - z_2)$
- z_2 对 a_1, w_2, b_2 的局部导数
- $\frac{dz_2}{da_1} = w_2, \frac{dz_2}{dw_2} = a_1, \frac{dz_2}{db_2} = 1$
- 因此 $\frac{dl}{da_1} = \frac{dl}{dz_2} \cdot w_2, \frac{dl}{dw_2} = \frac{dl}{dz_2} \cdot a_1, \frac{dz_2}{db_2} = \frac{dl}{dz_2}$



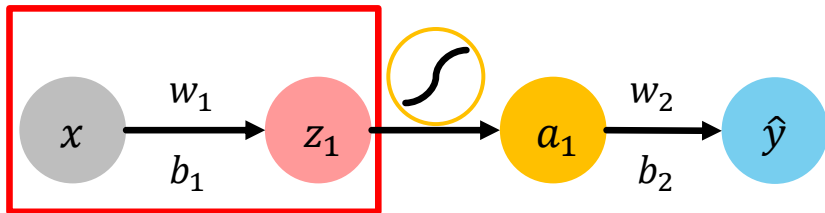
回到 神经网络

- a_1 的上游导数 $\frac{dl}{da_1}$ 已经在上一层计算好
- a_1 对 z_1 的局部导数
- $\frac{da_1}{dz_1} = \sigma'(z_1) = \sigma(z_1)[1 - \sigma(z_1)] = a_1(1 - a_1)$
- 因此 $\frac{dl}{dz_1} = \frac{dl}{da_1} \cdot a_1(1 - a_1)$



回到神经网络

- z_1 的上游导数 $\frac{dl}{dz_1}$ 已经在上一层计算好
- z_1 对 w_1, b_1 的局部导数
- $\frac{dz_1}{dw_1} = x, \frac{dz_1}{db_1} = 1$
- 因此 $\frac{dl}{dw_1} = \frac{dl}{dz_1} \cdot x, \frac{dl}{db_1} = \frac{dl}{dz_1}$



扩展

- 这一方法可以很容易扩展到任意层数的情形
- 但每层有多个神经元的情况怎么办？
- 此时 W_i 是矩阵, b_i 是向量

先上结论

规则0

- 参数导数的维度和参数的维度保持一致

\times

\times

规则1

- 对于逐元素计算的激活函数, $a = \sigma(z)$
- $a = (a^1, \dots, a^m)^T$, $z = (z^1, \dots, z^m)^T$, $a^i = \sigma(z^i)$
- 导数也可逐元素计算
- 假设 a 的上游导数为 $\frac{dl}{da} = \left(\frac{dl}{da^1}, \dots, \frac{dl}{da^m} \right)^T$, 则
- $\frac{dl}{dz} = \left(\frac{dl}{da^1} \cdot \sigma'(z^1), \dots, \frac{dl}{da^m} \cdot \sigma'(z^m) \right)^T = \frac{dl}{da} \circ \sigma'(z)$

简写

两个向量逐元素相乘

规则2

- 对于线性变换, $z = Wa + b$, $W \in \mathbb{R}^{n \times m}$
- $a = (a^1, \dots, a^m)^T$, $z = (z^1, \dots, z^n)^T$
- 假设 z 的上游导数为 $\frac{dl}{dz} = \left(\frac{dl}{dz^1}, \dots, \frac{dl}{dz^n} \right)^T$
- 那么
- $\frac{dl}{dW} = \frac{dl}{dz} a^T$, $\frac{dl}{db} = \frac{dl}{dz}$, $\frac{dl}{da} = W^T \frac{dl}{dz}$

规则2

- 对于线性变换, $z = Wa + b$, $W \in \mathbb{R}^{n \times m}$
- $a = (a^1, \dots, a^m)^T$, $z = (z^1, \dots, z^n)^T$
- 假设 z 的上游导数为 $\frac{dl}{dz} = \left(\frac{dl}{dz^1}, \dots, \frac{dl}{dz^n} \right)^T$
- 那么

$$\underbrace{\frac{dl}{dw}}_{n \times m} = \underbrace{\frac{dl}{dz}}_{n \times 1} \underbrace{a^T}_{1 \times m}, \quad \frac{dl}{db} = \frac{dl}{dz}, \quad \frac{dl}{da} = W^T \frac{dl}{dz}$$

规则2

- 对于线性变换, $z = Wa + b$, $W \in \mathbb{R}^{n \times m}$
- $a = (a^1, \dots, a^m)^T$, $z = (z^1, \dots, z^n)^T$
- 假设 z 的上游导数为 $\frac{dl}{dz} = \left(\frac{dl}{dz^1}, \dots, \frac{dl}{dz^n} \right)^T$
- 那么
- $\frac{dl}{dW} = \frac{dl}{dz} a^T$, $\boxed{\frac{dl}{db}} = \boxed{\frac{dl}{dz}}$, $\frac{dl}{da} = W^T \frac{dl}{dz}$
 $n \times 1$ $n \times 1$

规则2

- 对于线性变换, $z = Wa + b$, $W \in \mathbb{R}^{n \times m}$
- $a = (a^1, \dots, a^m)^T$, $z = (z^1, \dots, z^n)^T$
- 假设 z 的上游导数为 $\frac{dl}{dz} = \left(\frac{dl}{dz^1}, \dots, \frac{dl}{dz^n} \right)^T$
- 那么
- $\frac{dl}{dW} = \frac{dl}{dz} a^T$, $\frac{dl}{db} = \frac{dl}{dz}$, $\frac{dl}{da} = \underbrace{\frac{dl}{da}}_{m \times 1} = \underbrace{W^T}_{m \times n} \underbrace{\frac{dl}{dz}}_{n \times 1}$

规则3

- 不要照搬数学表达式
- 实际编程中要考虑到数据的存储方式



理论 VS 实际

规则0

理论

- 参数导数的维度和参数的维度保持一致

实际

- 参数导数的维度和参数的维度保持一致

规则1

理论

- 对于逐元素计算的激活函数，导数也可逐元素计算

- $a: [m \times 1], z: [m \times 1], a_i = \sigma(z_i)$

- $\frac{dl}{dz} = \frac{dl}{da} \circ \sigma'(z)$

实际

- 对于逐元素计算的激活函数，导数也可逐元素计算

- $A: [n \times m], Z: [n \times m], A_{ij} = \sigma(Z_{ij})$

- $\frac{dl}{dZ} = \frac{dl}{dA} \circ \sigma'(Z)$

规则1.5

$$\underset{2 \times 1}{1_n} \underset{1 \times 3}{b^T} = \begin{bmatrix} b_1 & b_2 & b_3 \\ b_1 & b_2 & b_3 \end{bmatrix}$$

加入数据的维度

实际

理论 W 的转置

左乘

★

→ funmax 图例

理论

左乘

■ 线性变换, $z = \boxed{W}a + b$, $W \in \mathbb{R}^{p \times m}$

■ $a: [m \times 1]$, $z: [p \times 1]$

■ $W: [p \times m]$, $b: [p \times 1]$

■ a 代表上一层神经元 (m 个)

■ z 代表下一层神经元 (p 个)

■ $Z = A\boxed{W} + \boxed{1_n}b^T$, $W \in \mathbb{R}^{m \times p}$

■ $A: [n \times m]$, $Z: [n \times p]$

■ $W: [m \times p]$, $b: [p \times 1]$

■ n 代表样本量

■ A 的第 i 行代表第 i 个观测在上一层的神经元

■ Z 的第 i 行代表第 i 个观测在下一层的神经元

规则2

理论

$$\blacksquare \frac{dl}{dw} = \frac{dl}{dz} a^T$$

$p \times m$ $p \times 1$ $1 \times m$

$$\blacksquare \frac{dl}{db} = \frac{dl}{dz}$$

$p \times 1$ $p \times 1$

$$\blacksquare \frac{dl}{da} = W^T \frac{dl}{dz}$$

$m \times 1$ $m \times p$ $p \times 1$

实际

$$\blacksquare \frac{dl}{dw} = A^T \frac{dl}{dz}$$

$m \times p$ $m \times n$ $n \times p$

$$\blacksquare \frac{dl}{db} = \left(\frac{dl}{dz} \right)^T \mathbf{1}_n$$

$p \times 1$ $p \times n$ $n \times 1$

$$\blacksquare \frac{dl}{dA} = \frac{dl}{dz} W^T$$

$n \times m$ $n \times p$ $p \times m$

自动微分 与优化

- 参见 [lec5-module.ipynb](#)

自动微分代码 $f(x, y) = x \cdot \log(x) + \sin(xy)$

```
import torch
```

```
x = torch.tensor([1.0], requires_grad=True)
```

矩阵、向量也可以

```
y = torch.tensor([2.0], requires_grad=True)
```

```
f = x * torch.log(x) + torch.sin(x * y)
```

```
f.backward()
```

backward 进行反向传播

```
x.grad
```

```
y.grad
```

不需要 pytorch 记录导数值时

```
with torch.no_grad():
```

```
    torch.log(x) + 1.0 + y * torch.cos(x * y)
```

问题: 行列式为正的矩阵 X , 定义 $f(X) = \log \det(X)$, 其中 \det 为 X 的行列式

那么 $\frac{df}{dX}$ 应该是什么?

$\frac{df}{dX}$ $\frac{d(\log \det(X))}{dX}$ $\frac{d(\log \det(X))}{dX}$

$$dX = \frac{dX}{dX} = 1$$

线性模型自动微分

$Y = X\beta + \epsilon$ 对 β 优化, 对 β 求导 参数!!!

$n \times p \quad p \times 1 \quad n \times 1$

损失函数: $\text{torch.mean}(\text{torch.square}(y - \text{yhat}))$

循环

$\text{nepoch} = 500$ 迭代次数

$\text{learning_rate} = 0.01$ 学习率

$\text{Losses} = []$ 损失值

for i in range(nepoch):

$\text{loss} = \text{loss_fn}(\text{bhat}, x, y)$ 这里获得 loss 为 tensor

$\text{loss.backward}()$ 求导获得 bhat.grad

$\text{losses.append}(\text{loss.item}())$ loss 具体值, 而不是 tensor
为了画图像

```

if i % 50 == 0:
    print(f"iteration {i}, loss = {loss.item()}, error = {torch.mean(
        torch.square(y - yhat))

with torch.no_grad():
    bhat -= learning_rate * bhat.grad

    # 清空梯度项
    bhat.grad = None

```

模块化编程

```
import torch.nn as nn
```

```
class MyModel(nn.Module)
```

```
    def __init__(self, beta_dim):
```

```
        super(MyModel, self).__init__()
```

```
        self.bhat = nn.Parameter(torch.zeros(beta_dim))
```

```
    def forward(self, x):
```

```
        yhat = torch.matmul(x, self.bhat)
```

参数, 自动微分

线性层很多层

预测值


```
return yhat
```

```
model = MyModel(beta_dim=p)  
print(list(model.parameters()))
```

```
nepoch = 500      迭代次数  
learning_rate = 0.01  学习率  
losses = []       损失值
```

```
opt = torch.optim.SGD(model.parameters(), lr=learning_rate)
```

对参数优化

```
for i in range(nepoch):
```

```
    yhat = model(x)  调用forward返回
```

```
    loss = torch.mean(torch.square(y - yhat))  
    losses.append(loss.item())
```

预测值

loss = loss_fn(bhat, x, y)

损失函数 与分布有关

```
    opt.zero_grad()
```

```
    loss.backward()  求导
```

with torch.no_grad():

bhat -= learning_rate * bhat.grad
清空梯度项

`opt.step()` 实现更新beta —

`bhat.grad = None`

`if i % 50 == 0:`

`print(f"iteration {i}, loss = {loss.item()} ,`

`h > 0`