

习题 4-1

我们有:

$$z = w^T X + b$$

$$a = \sigma(z)$$

计算梯度:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w} &= \frac{\partial z}{\partial w} \frac{\partial a}{\partial z} \frac{\partial \mathcal{L}}{\partial a} \\ &= X^T \cdot \sigma'(z) \cdot \frac{\partial \mathcal{L}}{\partial a} \end{aligned}$$

考虑到:

$$x_i > 0, \forall x_i \in X$$

$$\sigma'(z) > 0$$

故:

$$\text{sign}\left(\frac{\partial \mathcal{L}}{\partial w}\right) = \text{sign}\left(\frac{\partial \mathcal{L}}{\partial a}\right)$$

即, 一次更新中, w 中 w_1, \dots, w_i 的正负号必然相同.

这在大部分情况下会导致参数更新的速率下降.

具体情况如图1所示, 引用自: [cs231n lecture7 slides](#)

习题 4-2

定义 XOR 问题如下:

输入:

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$x_i \in \{0, 1\}$$

输出:

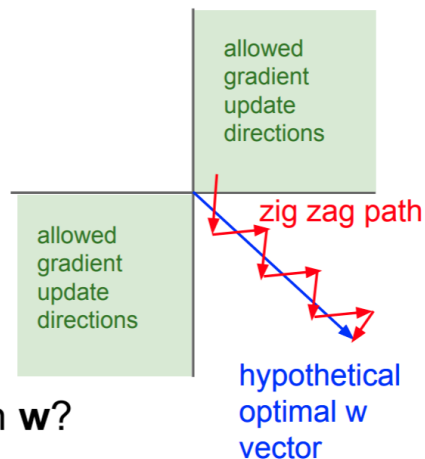
$$\hat{y} = \begin{cases} 1 & \text{if } x_1 \neq x_2 \\ 0 & \text{otherwise} \end{cases}$$

Consider what happens when the input to a neuron is always positive...

$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on \mathbf{w} ?

Always all positive or all negative :(



Fei-Fei Li, Ranjay Krishna, Danfei Xu

Lecture 7 - 34

April 20, 2021

图 1: X 始终同号情况下的梯度更新情况

根据题意, 我们拥有以下模型参数:

$$w^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \end{bmatrix}$$

$$w^{(2)} = \begin{bmatrix} w_1^{(2)} & w_2^{(2)} \end{bmatrix}$$

$$b^{(2)} = \begin{bmatrix} b^{(2)} \end{bmatrix}$$

故:

$$z^{(1)} = w^{(1)}x + b^{(1)}$$

$$a^{(1)} = \text{relu}(z^{(1)})$$

$$\hat{y} = w^{(2)}a^{(1)} + b^{(2)}$$

则其中一解为:

$$\begin{aligned}w^{(1)} &= \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \\b^{(1)} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\w^{(2)} &= \begin{bmatrix} 1 & 1 \end{bmatrix} \\b^{(2)} &= \begin{bmatrix} 0 \end{bmatrix}\end{aligned}$$

参数为以上解时, 则:

x	$z^{(1)}$	$a^{(1)}$	\hat{y}
$[0, 0]^T$	$[0, 0]^T$	$[0, 0]^T$	0
$[0, 1]^T$	$[-1, 1]^T$	$[0, 1]^T$	1
$[1, 0]^T$	$[1, -1]^T$	$[1, 0]^T$	1
$[1, 1]^T$	$[0, 0]^T$	$[0, 0]^T$	0

习题 4-3

对第一个隐藏层的第一个神经元, 我们有:

$$\begin{aligned}z_1^{(1)} &= w_1^{(1)}x + b_1^{(1)} \\a_1^{(1)} &= \text{relu}(z_1^{(1)}) \\&= \max(0, z_1^{(1)})\end{aligned}$$

现假设, 对训练集中任意 x , 都有 $z_1^{(1)} \leq 0$.

考虑更新参数时计算的梯度:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_1^{(1)}} &= \frac{\partial z_1^{(1)}}{\partial w_1^{(1)}} \frac{\partial a_1^{(1)}}{\partial z_1^{(1)}} \frac{\partial \mathcal{L}}{\partial a_1^{(1)}} \\&= x^T \cdot \mathbf{1}\{z_1^{(1)} > 0\} \cdot \frac{\partial \mathcal{L}}{\partial a_1^{(1)}} \\&= x^T \cdot 0 \cdot \frac{\partial \mathcal{L}}{\partial a_1^{(1)}} \\&= 0\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial b_1^{(1)}} &= \frac{\partial z_1^{(1)}}{\partial b_1^{(1)}} \frac{\partial a_1^{(1)}}{\partial z_1^{(1)}} \frac{\partial \mathcal{L}}{\partial a_1^{(1)}} \\
&= 1 \cdot \mathbf{1}\{z_1^{(1)} > 0\} \cdot \frac{\partial \mathcal{L}}{\partial a_1^{(1)}} \\
&= 1 \cdot 0 \cdot \frac{\partial \mathcal{L}}{\partial a_1^{(1)}} \\
&= 0
\end{aligned}$$

考虑更新参数时的情况:

$$\begin{aligned}
w_1^{(1)} &:= w_1^{(1)} - \alpha \frac{\partial \mathcal{L}}{\partial w_1^{(1)}} \\
&= w_1^{(1)} \\
b_1^{(1)} &:= b_1^{(1)} - \alpha \frac{\partial \mathcal{L}}{\partial b_1^{(1)}} \\
&= b_1^{(1)}
\end{aligned}$$

这意味着, 在该数据集中, $w_1^{(1)}$ 和 $b_1^{(1)}$ 将永远不会被更新, 且对整个模型毫无作用. 即, 它死亡了.

一种解决方案为: 替换 rule 函数或对其进行修改, 使 $z_1^{(1)} \leq 0$ 时, 通过激活函数依旧能输出一个微小的正数, 以保证进行反向传播时, 其梯度不为 0.

习题 4-5

令参数数量为 n :

$$n = \frac{N}{L}(M_0 + 1) + (L - 1)\frac{N}{L}\left(\frac{N}{L} + 1\right) + 1 \cdot \left(\frac{N}{L} + 1\right)$$

其中:

- $\frac{N}{L}(M_0 + 1)$ 为输入层到第一层隐藏层之间的参数
- $(L - 1)\frac{N}{L}\left(\frac{N}{L} + 1\right)$ 为第一层隐藏层到最后一层隐藏层之间的参数
- $1 \cdot \left(\frac{N}{L} + 1\right)$ 为最后一层隐藏层到输出层间的参数

习题 4-7

可以, 但这样做并不一定带来相应的好处.

当我们使用正则化时, 通常是为了避免过拟合. 考虑 $y = wx + b$, 我们可以发现过拟合主要是由 w 对特定数据过于敏感导致的, 而在这一问题中, b 并未扮演重要的角色 (w 是高维向量, 而 b 只是一个单一数值). 故对 b 进行正则化处理并不能造成多少效果.

引用自: [Coursera - Improving Deep Neural Networks - Regularization](#)

习题 4-8

let

$$\begin{aligned} z^{[l]} &= w^{[l]}a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g(z^{[l]}) \\ z^{[l+1]} &= w^{[l+1]}a^{[l]} + b^{[l+1]} \\ a^{[l+1]} &= g(z^{[l+1]}) \end{aligned}$$

then

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w^{[l]}} &= \frac{\partial z^{[l]}}{\partial w^{[l]}} \frac{\partial a^{[l]}}{\partial z^{[l]}} \frac{\partial z^{[l+1]}}{\partial a^{[l]}} \frac{\partial a^{[l+1]}}{\partial z^{[l+1]}} \frac{\partial \mathcal{L}}{\partial a^{[l+1]}} \\ &= a^{[l-1]} g'(z^{[l]}) w^{[l+1]} g'(z^{[l+1]}) \frac{\partial \mathcal{L}}{\partial a^{[l+1]}} \end{aligned}$$

if $w_i^{[l+1]} = 0, \forall w_i^{[l+1]} \in w^{[l+1]}$, then $\frac{\partial \mathcal{L}}{\partial w^{[l]}} = 0$. This means that parameters will never be updated.

[Initializing neural networks](#) This website has a good online visualisation tool.

习题 4-9

Here has some disscussion: [Why can't we handle vanishing gradient problem in neural nets using large step sizes?](#)