



数值稳定性





神经网络的梯度

- 考虑如下有 d 层的神经网络

$$\mathbf{h}^t = f_t(\mathbf{h}^{t-1}) \quad \text{and} \quad y = \ell \circ f_d \circ \dots \circ f_1(\mathbf{x})$$

- 计算损失 ℓ 关于参数 \mathbf{W}_t 的梯度

$$\frac{\partial \ell}{\partial \mathbf{W}^t} = \frac{\partial \ell}{\partial \mathbf{h}^d} \underbrace{\frac{\partial \mathbf{h}^d}{\partial \mathbf{h}^{d-1}} \cdots \frac{\partial \mathbf{h}^{t+1}}{\partial \mathbf{h}^t}}_{\text{d-t 次矩阵乘法}} \frac{\partial \mathbf{h}^t}{\partial \mathbf{W}^t}$$

d-t 次矩阵乘法

数值稳定性的常见两个问题

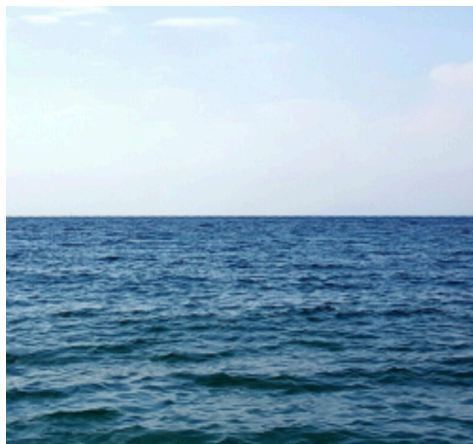
$$\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i}$$

梯度爆炸



$$1.5^{100} \approx 4 \times 10^{17}$$

梯度消失



$$0.8^{100} \approx 2 \times 10^{-10}$$



例子：MLP

- 加入如下 MLP （为了简单省略了偏移）

$$f_t(\mathbf{h}^{t-1}) = \sigma(\mathbf{W}^t \mathbf{h}^{t-1}) \quad \sigma \text{ 是激活函数}$$

$$\frac{\partial \mathbf{h}^t}{\partial \mathbf{h}^{t-1}} = \text{diag}(\sigma'(\mathbf{W}^t \mathbf{h}^{t-1}))(\mathbf{W}^t)^T \quad \sigma' \text{ 是 } \sigma \text{ 的导数函数}$$

$$\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i} = \prod_{i=t}^{d-1} \text{diag}(\sigma'(\mathbf{W}^i \mathbf{h}^{i-1}))(\mathbf{W}^i)^T$$

梯度爆炸



- 使用 ReLU 作为激活函数

$$\sigma(x) = \max(0, x) \quad \text{and} \quad \sigma'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

- $\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i} = \prod_{i=t}^{d-1} \text{diag}(\sigma'(\mathbf{W}^i \mathbf{h}^{i-1})) (\mathbf{W}^i)^T$ 的一些元素会来自于 $\prod_{i=t}^{d-1} (\mathbf{W}^i)^T$
 - 如果 $d-t$ 很大, 值将会很大



梯度爆炸的问题

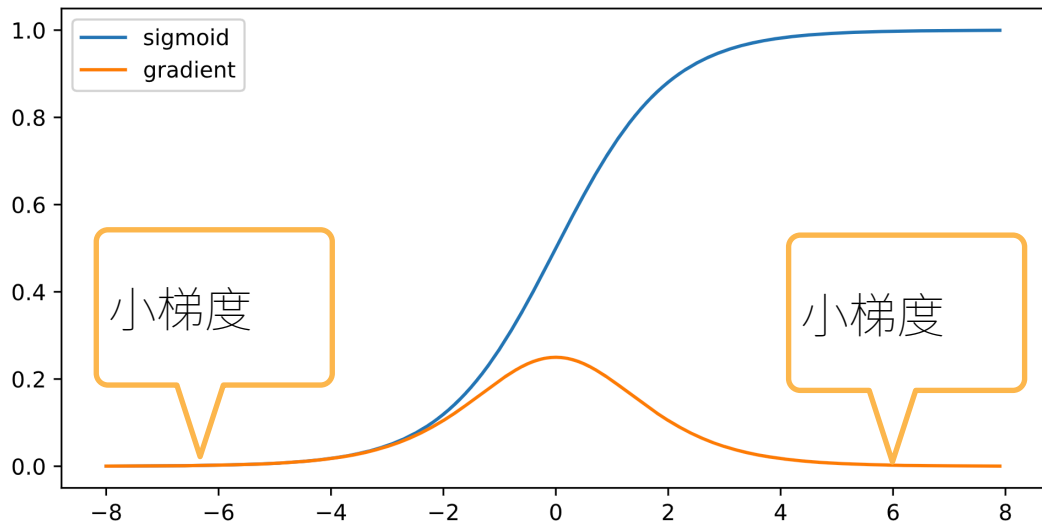
- 值超出值域 (infinity)
 - 对于 16位浮点数尤为严重 (数值区间 $6e-5$ - $6e4$)
- 对学习率敏感
 - 如果学习率太大 -> 大参数值 -> 更大的梯度
 - 如果学习率太小 -> 训练无进展
 - 我们可能需要在训练过程不断调整学习率

梯度消失



- 使用 sigmoid 作为激活函数

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$



梯度消失



- 使用 sigmoid 作为激活函数

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

- $\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i} = \prod_{i=t}^{d-1} \text{diag}(\sigma'(\mathbf{W}^i \mathbf{h}^{i-1}))(W^i)^T$ 的元素值是 d-t 个小数值的乘积

$$0.8^{100} \approx 2 \times 10^{-10}$$



梯度消失的问题

- 梯度值变成 0
 - 对 16 位浮点数尤为严重
- 训练没有进展
 - 不管如何选择学习率
- 对于底部层尤为严重
 - 仅仅顶部层训练的较好
 - 无法让神经网络更深

总结



- 当数值过大或者过小时会导致数值问题
- 常发生在深度模型中，因为其会对 n 个数累乘