



deeplearning.ai

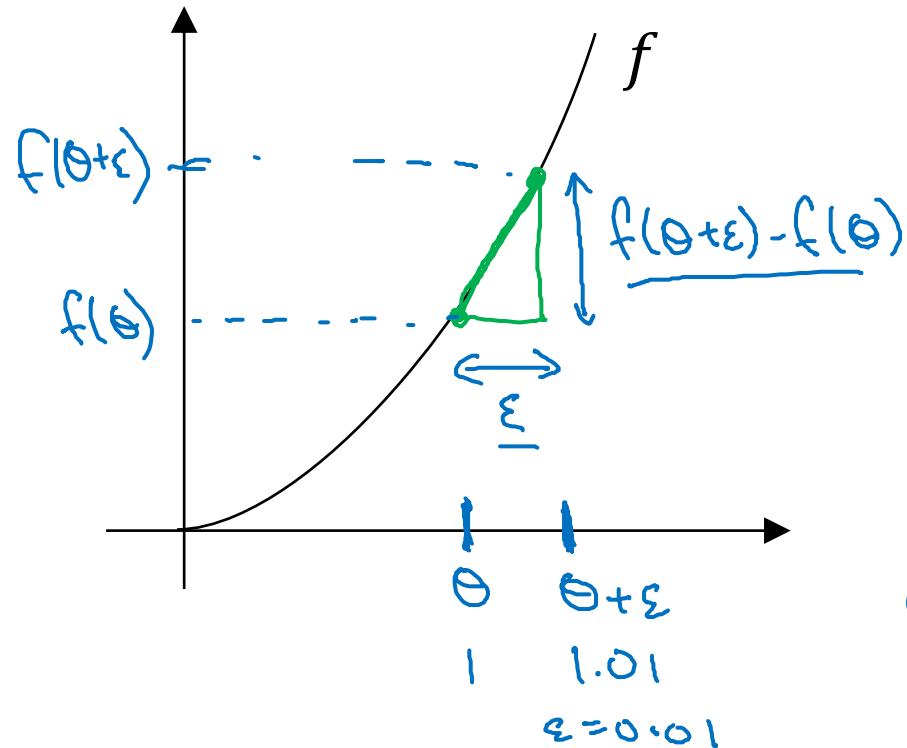
# Setting up your optimization problem

---

## Numerical approximation of gradients

# Checking your derivative computation

I  $f(\theta) = \theta^3$   
 $\theta \in \mathbb{R}.$



$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

$\frac{dw}{db}$   $\rightarrow$   $g(\theta) = 3\theta^2$

$g(\theta) = 3 \cdot (1)^2 = 3$   
 when  $\theta = 1$

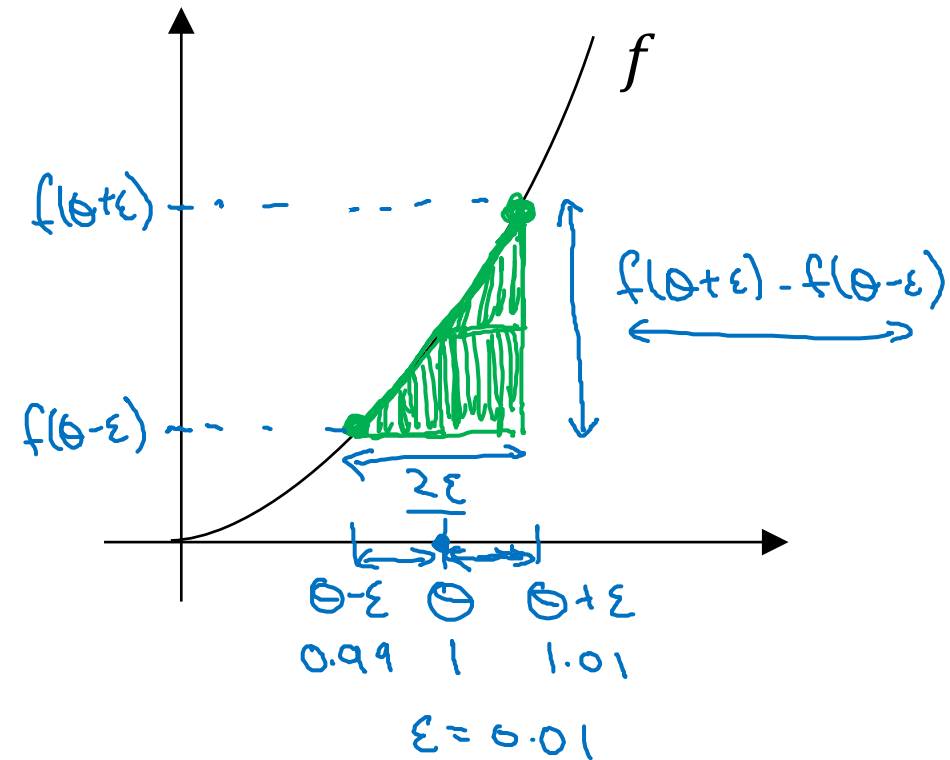
$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - 1^3}{0.01} = \frac{3.0301}{0.01} = 3.0301 \approx 3$$

Annotations:  $\theta = 1$ ,  $\theta + \epsilon = 1.01$ ,  $\epsilon = 0.01$ . The calculation shows  $3.1$  and  $3.2$  as intermediate steps.

# Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$



$$\left[ \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx \underline{g(\theta)} \right]$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

---


$$\left\{ \begin{array}{l} f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \quad \begin{array}{l} \mathcal{O}(\epsilon^2) \\ 0.01 \\ \underline{0.0001} \end{array} \quad \left| \quad \begin{array}{l} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \quad \text{error: } \mathcal{O}(\epsilon) \\ \uparrow \quad \uparrow \\ 0.01 \end{array} \end{array} \right.$$

当你实现梯度逆传播时 你会发现一个测试 名叫梯度检验 可以帮助你确保 你的梯度逆传播的实现是正确的 因为有时写下这些方程后 你无法百分百确定 你的所有细节都做对了 实现了正确的逆传播 所以 为了之后能介绍梯度检验 我们先谈谈如何在数值上近似计算梯度 在下一个视频中 我们再讨论如何实现 梯度检验 从而确保逆传播的实现是正确的 我们来看一下函数  $f$  我把它画在了这里 这个  $f(\theta)$  等于  $\theta$  立方 我们从某些  $\theta$  取值出发 比如  $\theta=1$  这次我们不再将  $\theta$  向右扰动到  $\theta+\varepsilon$  我们将同时向左右扰动 得到  $\theta-\varepsilon$  和  $\theta+\varepsilon$  所以这里是1 这里是1.01 这里是0.99  $\varepsilon$  与之前一样是0.01 事实上 相比于取这个小三角形 并计算高和宽 你还能得到更好的估计值 如果你选取这个点  $f(\theta-\varepsilon)$  和这个点 然后计算这个大三角形的高和宽 由于一些技术原因 我在此不作解释 这个大三角形的 高和宽会给你一个  $\theta$  点处梯度更好的近似值 很容易发现 如果不取右上的这个小三角形 你可以认为你有两个小三角形 右上的这个和左下的这个 你相当于同时考虑这两个三角形 通过使用这个大的绿三角形 所以你这次取了双侧的差值而不是单侧的差值 我们来写成数学形式 这个点是  $f(\theta+\varepsilon)$  而这个点是  $f(\theta-\varepsilon)$  所以这个大的绿三角形的高是  $f(\theta+\varepsilon)-f(\theta-\varepsilon)$  所以这个大的绿三角形的高是  $f(\theta+\varepsilon)-f(\theta-\varepsilon)$  至于宽度 这里是1个  $\varepsilon$  这里是第2个  $\varepsilon$  所以这个绿色三角形的宽就是  $2\varepsilon$  所以这里的高度就是 首先是高度 也就是  $f(\theta+\varepsilon)-f(\theta-\varepsilon)$  除以宽度 也就是  $2\varepsilon$  我写在这里了 这个很可能与  $g(\theta)$  很接近 代入这些值 记住  $f(\theta)$  是  $\theta$  立方 这里  $\theta+\varepsilon$  也就是1.01 我取它的立方 然后减去0.99的立方 再除以  $2*0.01$  你可以暂停视频并在计算器上算一下 你会得到结果是3.0001 前一个幻灯片中 我们知道  $g(\theta)$  也就是  $3\theta^2$  而  $\theta$  是1 所以这两个值十分接近 近似误差是0.0001 在上一个幻灯片中 我们取过单侧的差值  $\theta$  和  $\theta+\varepsilon$  之间的结果是3.0301 所以近似误差就是0.03而不是0.0001 所以用这个取双侧差值的方法 来近似导数 你会发现结果非常接近3 所以这将让你更加自信  $g(\theta)$  很可能就是求  $f$  导数的正确实现 当你把这个方法用于梯度检验和逆传播时 它运行起来很可能比用单侧差值要慢两倍 但从实践的角度 我认为这个方法值得一用 因为它精确很多 我再说一些选修的理论知识 给你们之中比较熟悉微积分的人 如果你听不懂我即将说到的内容也没有关系 事实上 导数的正式定义 就是对于很小的  $\varepsilon$  计算  $[f(\theta+\varepsilon)-f(\theta-\varepsilon)]/(2\varepsilon)$  就是对于很小的  $\varepsilon$  计算  $[f(\theta+\varepsilon)-f(\theta-\varepsilon)]/(2\varepsilon)$  而导数的正式定义就是右边这个公式 当  $\varepsilon$  趋近于0时的极限 极限的定义就是微积分课上学的那样 但我这里不再详述 对于一个非零的  $\varepsilon$  值 你可以证明这个近似的误差 在  $\varepsilon$  平方这个阶上  $\varepsilon$  是个很小的数 如果  $\varepsilon$  是0.01 就像这里 那么  $\varepsilon$  平方就是0.0001 这个大O记号就表示误差就是某个常数乘以这个 这就是我们的近似误差 这个例子中的大O的常数恰好就是1 相比而言 如果我们用这边的另一个公式 误差就在  $\varepsilon$  这个阶上 当  $\varepsilon$  是一个小于1的数时  $\varepsilon$  就比  $\varepsilon$  平方大很多 这也就是为什么 这个公式不如左边这个公式精确 这也就是为什么我们做梯度检验时采用双侧差值 你计算  $f(\theta+\varepsilon)-f(\theta-\varepsilon)$  再除以  $2\varepsilon$  而不使用这个不够精确的单侧差值 如果你不理解我最后说的两点 所有东西都写在这儿了 Don't worry about it. Don't worry about it. 对微积分和数值近似比较熟悉的人 可以多学一些 简单来说就是双侧差值的公式更加精确 我们在下一个梯度检验视频中就会用到这个 所以你们学习了如何取双侧差值 来在数值上验证是否给定的函数  $g(\theta)$  是函数  $f$  的导数的正确实现 我们来看看如何使用这个来验证是否 你的逆传播的实现是正确的 还是说里面有错误要剔除掉