

3 导数、反向传播和复杂度

概要

- 矩阵微积分
- 链式法则
- 自动微分法
- 反向传播

矩阵微积分

$\frac{d}{dx} \left[\frac{f(x)}{g(x)} \right] = \frac{g(x)f'(x) - f(x)g'(x)}{g(x)^2}$
 $F = mg = ma = m \frac{d^2h}{dt^2}$
 $\frac{dA}{dt} = \frac{dB}{dt} = -\frac{dC}{dt} = -\frac{dD}{dt} = (c_1)T^{\frac{1}{2}}AB - (c_2)T^{\frac{1}{2}}CD$
 $\frac{dA}{dt} = \frac{dB}{dt} = \frac{dC}{dt} = \frac{dD}{dt} = (c_1)AB - (c_2)CD$
 $y = mx + b$
 $\frac{d^2x}{dt^2} = -kx$

$\frac{dx}{dx} = \frac{dy}{dy} = \frac{dz}{dz}$
 $\int \frac{1}{x} dx = \ln|x| + C$
 $\int \sin x dx = -\cos x + C$
 $\int_a^b f'(x) dx = f(b) - f(a)$
 $\frac{d^2x}{dt^2} = -kx$

Calculus

$x^2 - 3x - 4 = 0$
 $4x^2 - 3x - 1 = 0$

$\int f(x) dx$

$\frac{dA}{dt} = \frac{dB}{dt} = -\frac{dC}{dt} = -\frac{dD}{dt} = (c_1)T^{\frac{1}{2}}AB - (c_2)T^{\frac{1}{2}}CD$

$x^2 = A \frac{dT}{dt} = (c_1) \frac{dA}{dt} - (c_2)(T_0 - T)$

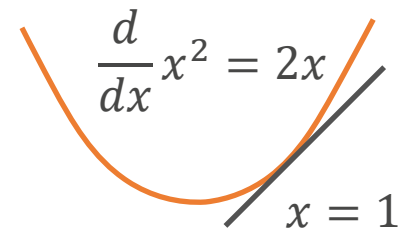
$\left[x + \frac{b}{2a} \right]^2 = \frac{b^2 - 4ac}{4a^2}$
 $x + \frac{b}{2a} = \frac{\sqrt{b^2 - 4ac}}{2a}$ or $x + \frac{b}{2a} = -\frac{\sqrt{b^2 - 4ac}}{2a}$
 $(x+h, f(x+h))$

$\frac{d}{dx} \int_a^x f(t) dt = f(x)$
 $\frac{d^2x}{dt^2} = -kx - f \frac{dx}{dt} + A \sin(\omega t)$
 $y' = v$, and $v' = -ky - fv + A \sin(\omega t)$
 $f(x-h) - f(x)$

标量求导回顾

y	a	x^n	$\exp(x)$	$\log(x)$	$\sin(x)$
$\frac{dy}{dx}$	0	nx^{n-1}	$\exp(x)$	$\frac{1}{x}$	$\cos(x)$

导数是切线的斜率



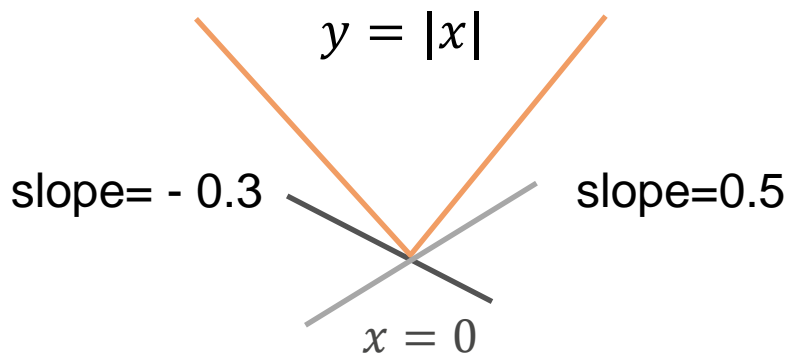
y	$u + v$	uv	$y = f(u), u = g(x)$
$\frac{dy}{dx}$	$\frac{du}{dx} + \frac{dv}{dx}$	$\frac{du}{dx}v + \frac{dv}{dx}u$	$\frac{dy}{du} \frac{du}{dx}$

切线的斜率为2

次导数

不可求导情况下的导数

Example 1:



$$\frac{\partial}{\partial x} \max(x, 0) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \\ a & \text{if } x = 0, a \in [0,1] \end{cases}$$

$$\frac{\partial |x|}{\partial x} = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \\ a & \text{if } x = 0, a \in [-1,1] \end{cases}$$

梯度

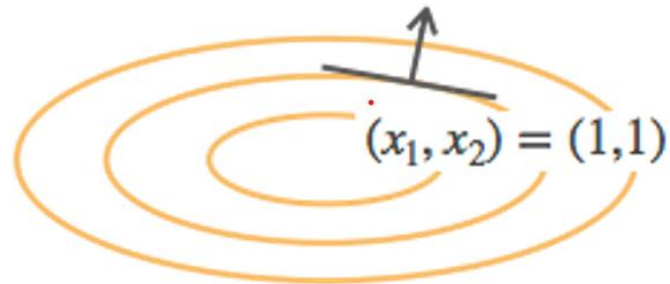
► 矢量求导推广

		标量	矢量
		x	\mathbf{x}
标量	y	$\frac{\partial y}{\partial x}$	$\frac{\partial y}{\partial \mathbf{x}}$
矢量	\mathbf{y}	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$

$\partial y / \partial \mathbf{x}$

$$\text{▶ } \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n} \right]$$

$$\text{▶ } \frac{\partial}{\partial \mathbf{x}} (x_1^2 + 2x_2^2) = [2x_1, 4x_2]$$



例子

y	a	au	$sum(\mathbf{x})$	$\ \mathbf{x}\ ^2$
$\frac{\partial y}{\partial \mathbf{x}}$	$\mathbf{0}^T$	$a \frac{\partial u}{\partial \mathbf{x}}$	$\mathbf{1}^T$	$2\mathbf{x}^T$

y	$u + v$	uv	$\langle \mathbf{u}, \mathbf{v} \rangle$
$\frac{\partial y}{\partial \mathbf{x}}$	$\frac{\partial u}{\partial \mathbf{x}} + \frac{\partial v}{\partial \mathbf{x}}$	$\frac{\partial u}{\partial \mathbf{x}} v + \frac{\partial v}{\partial \mathbf{x}} u$	$\mathbf{u}^T \frac{\partial \mathbf{v}}{\partial \mathbf{x}} + \mathbf{v}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$

$\partial \mathbf{y} / \partial x$

$$\triangleright \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}, \quad \frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}$$

- ▶ $\partial \mathbf{y} / \partial x$ 分子布局 (numerator-layout 或 Jacobian formulation), 是行矢量
- ▶ $\partial \mathbf{y} / \partial \mathbf{x}$ 分母布局 (denominator-layout 或 Hessian formulation), 是列矢量

$$\triangleright \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{y} \in \mathbb{R}^m, \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{m \times n}$$

$$\triangleright \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \frac{\partial y_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1}, \frac{\partial y_1}{\partial x_2}, \dots, \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1}, \frac{\partial y_2}{\partial x_2}, \dots, \frac{\partial y_2}{\partial x_n} \\ \vdots \\ \frac{\partial y_m}{\partial x_1}, \frac{\partial y_m}{\partial x_2}, \dots, \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

例子

▶ $\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m, \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{m \times n}$

▶ \mathbf{y}, \mathbf{a} 和 \mathbf{A} 不是关于 \mathbf{x} 的系数 $\mathbf{x}^T \mathbf{A}$

▶ $\mathbf{0}$ 和 \mathbf{I} 为矩阵

$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	$\mathbf{0}$	\mathbf{I}	\mathbf{A}	\mathbf{A}^T
---	--------------	--------------	--------------	----------------

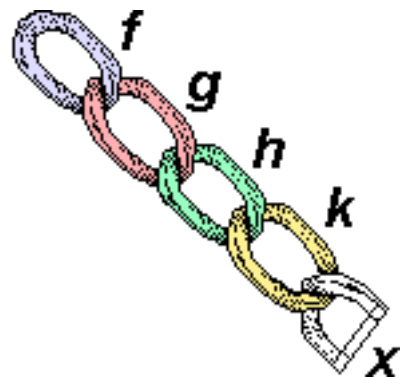
\mathbf{y}	$a\mathbf{u}$	$\mathbf{A}\mathbf{u}$	$\mathbf{u} + \mathbf{v}$
--------------	---------------	------------------------	---------------------------

$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	$a \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$	$\mathbf{A} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \frac{\partial \mathbf{v}}{\partial \mathbf{x}}$
---	---	--	---

推广到矩阵

		标量	矢量	矩阵
		$x(1,)$	$\mathbf{x}(n, 1)$	$\mathbf{X}(n, k)$
标量	$y(1,)$	$\frac{\partial y}{\partial x}(1,)$	$\frac{\partial y}{\partial \mathbf{x}}(1, n)$	$\frac{\partial y}{\partial \mathbf{X}}(k, n)$
矢量	$\mathbf{y}(m, 1)$	$\frac{\partial \mathbf{y}}{\partial x}(m, 1)$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}(m, n)$	$\frac{\partial \mathbf{y}}{\partial \mathbf{X}}(m, k, n)$
矩阵	$\mathbf{Y}(m, l)$	$\frac{\partial \mathbf{Y}}{\partial x}(m, l)$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}(m, l, n)$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}(m, l, k, n)$

链式法则



链式法则

➤ 链式法则 – 标量:

$$y = f(u), u = g(x), \frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x}$$

➤ 链式法则 – 矢量:

$$\begin{array}{ccc} \frac{\partial y}{\partial \mathbf{x}} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial \mathbf{x}} & \frac{\partial y}{\partial \mathbf{x}} = \frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} & \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \\ (1, n)(1, 1)(1, n) & (1, n)(1, k)(k, n) & (m, n)(m, k)(k, n) \end{array}$$

例1

▶ 假设 $\mathbf{x}, \mathbf{w} \in \mathbb{R}^n, y \in \mathbb{R} z = (\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$

计算 $\frac{\partial z}{\partial \mathbf{w}}$

例1

➤ 假设 $\mathbf{x}, \mathbf{w} \in \mathbb{R}^n, y \in \mathbb{R} z = (\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$

➤ 计算 $\frac{\partial z}{\partial \mathbf{w}}$

➤ 分解

➤ $a = \langle \mathbf{x}, \mathbf{w} \rangle$

➤ $b = a - y$

➤ $z = b^2$

➤ 求偏导

$$\frac{\partial z}{\partial \mathbf{w}} = \frac{\partial z}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial \mathbf{w}} = \frac{\partial b^2}{\partial b} \frac{\partial a - y}{\partial a} \frac{\partial \langle \mathbf{x}, \mathbf{w} \rangle}{\partial \mathbf{w}} = 2b \cdot 1 \cdot \mathbf{x}^T = 2(\langle \mathbf{x}, \mathbf{w} \rangle - y)\mathbf{x}^T$$

例2

假设 $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{w} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$ $z = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

计算 $\frac{\partial z}{\partial \mathbf{w}}$

分解 $\mathbf{a} = \mathbf{X}\mathbf{w}$

$$\mathbf{b} = \mathbf{a} - \mathbf{y}$$

$$z = \|\mathbf{b}\|^2$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$$

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{w}} &= \frac{\partial z}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{w}} \\ &= \frac{\partial \|\mathbf{b}\|^2}{\partial \mathbf{b}} \frac{\partial \mathbf{a} - \mathbf{y}}{\partial \mathbf{a}} \frac{\partial \mathbf{X}\mathbf{w}}{\partial \mathbf{w}} \\ &= 2\mathbf{b}^T \times \mathbf{I} \times \mathbf{X} \\ &= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T \mathbf{X} \end{aligned}$$

自动微分法



自动微分 (AD)

➤ 自动微分(AD)将符号微分法应用于最基本的算子，然后代入数值，应用于整个函数

➤ 其它常见微分法

➤ 符号微分法

➤ In[1] := D[4x³ + x² + 3, x]

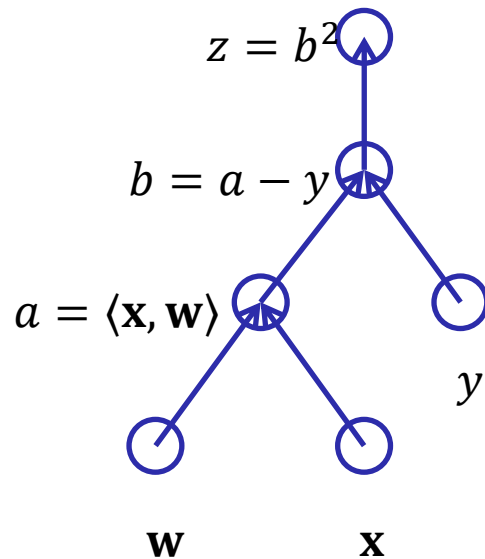
➤ Out[1] = 2x + 12x²

➤ 数值微分法

$$\text{➤ } \frac{\partial f(x)}{\partial x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

计算图

- 分解成最基本的方程
- 构造有向无环图来表示运算
- 假设 $z = (\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$



计算图

➤ 分解成最基本的方程

➤ 构造有向无环图来表示运算

➤ 显性构造Tensorflow/Theano/MXNet

➤ 隐性构造

PyTorch/MXNet

```
from mxnet import autograd, nd
```

```
with autograd.record():
```

```
    a = nd.ones((2,1))
```

```
    b = nd.ones((2,1))
```

```
    c = 2 * a + b
```

两种模式

➤ 链式法则

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u_n} \frac{\partial u_n}{\partial u_{n-1}} \cdots \frac{\partial u_2}{\partial u_1} \frac{\partial u_1}{\partial x}$$

➤ 正向传播

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u_n} \left(\frac{\partial u_n}{\partial u_{n-1}} \left(\cdots \left(\frac{\partial u_2}{\partial u_1} \frac{\partial u_1}{\partial x} \right) \right) \right)$$

➤ 【每次计算都是针对 x 求偏导，扩展求偏导函数】

➤ 与函数求值得顺序相同，求导所需函数值可以和对应的微分一起求出来

➤ 反向传播 【固定求偏导函数，扩展偏导变量，手工求导方式】

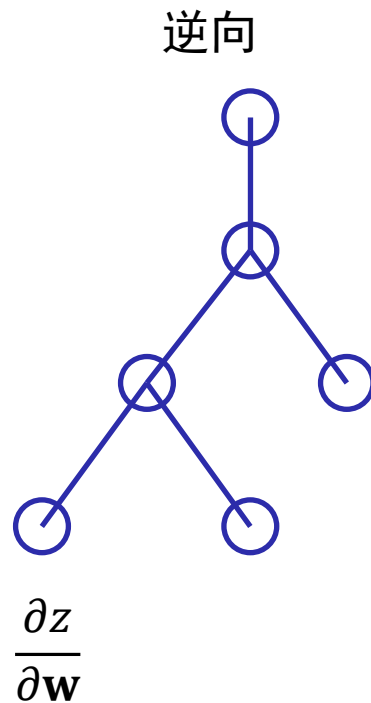
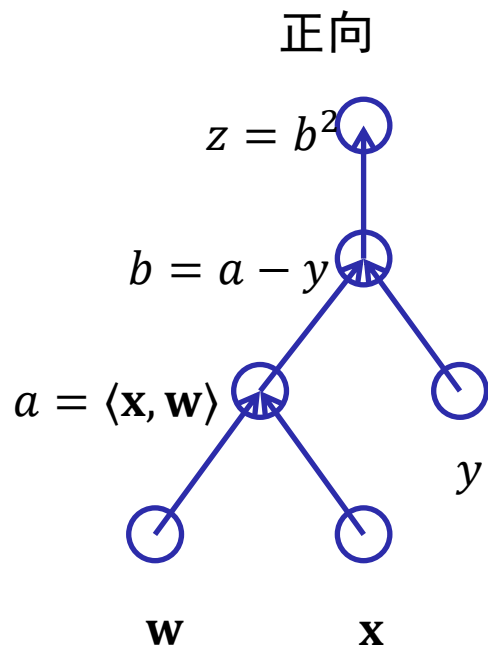
$$\frac{\partial y}{\partial x} = \left(\left(\left(\frac{\partial y}{\partial u_n} \frac{\partial u_n}{\partial u_{n-1}} \right) \cdots \right) \frac{\partial u_2}{\partial u_1} \right) \frac{\partial u_1}{\partial x}$$

正向反向的比较

- 求值必须正向，求导有两种做法。
- 正向：
 - 从 w 和 x 出发，由下往上求偏导，直到 z
 - 每次前向计算只能计算对一个自变量的偏导数。如果有一个函数，其输入有 n 个，求解整个函数梯度需要 n 遍计算过程
- 反向
 - 从 z 出发，自顶向下，求到该节点的变量的偏导
 - 子节点可以使用父节点的计算结果。且，某些分支，如 $\partial z / \partial y$ 可以忽略
 - 反向模式的优点：
 - 通过一次反向传输，就计算出所有偏导数，中间的偏导数计算只需计算一次，减少了重复计算的工作量，在多参数的时候后向自动微分的时间复杂度更低
 - 反向模式的缺点：

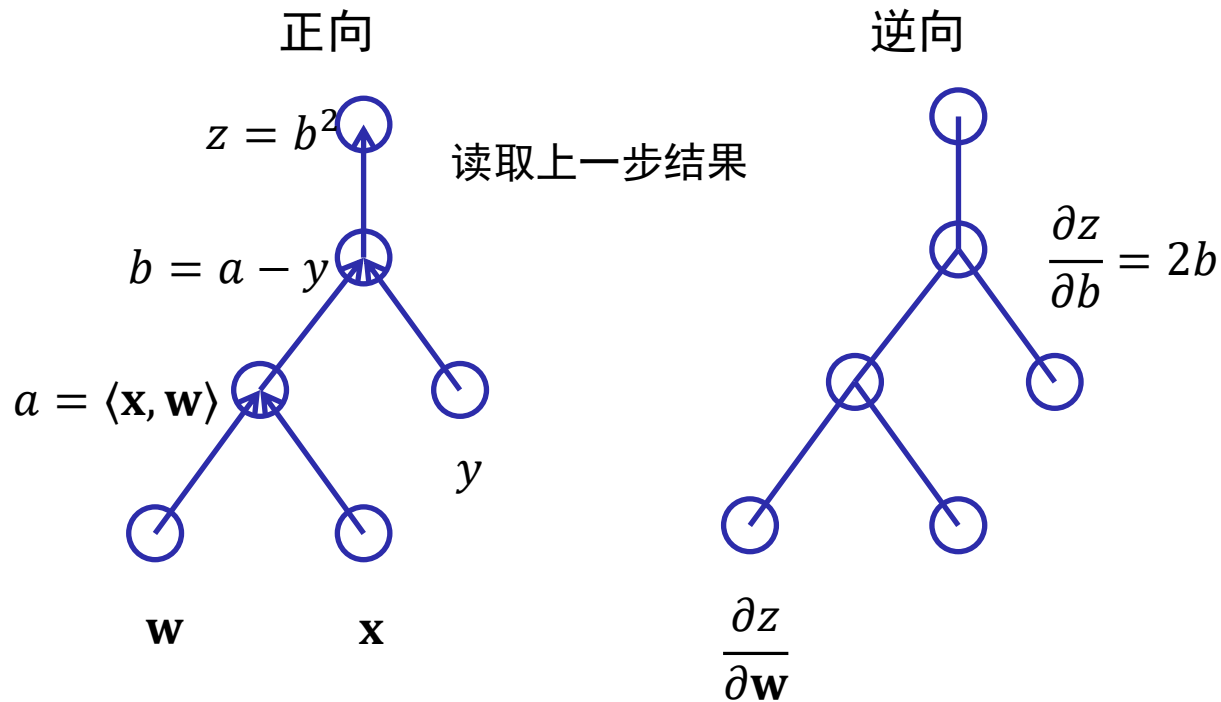
反向传播

► 假设 = $(\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$



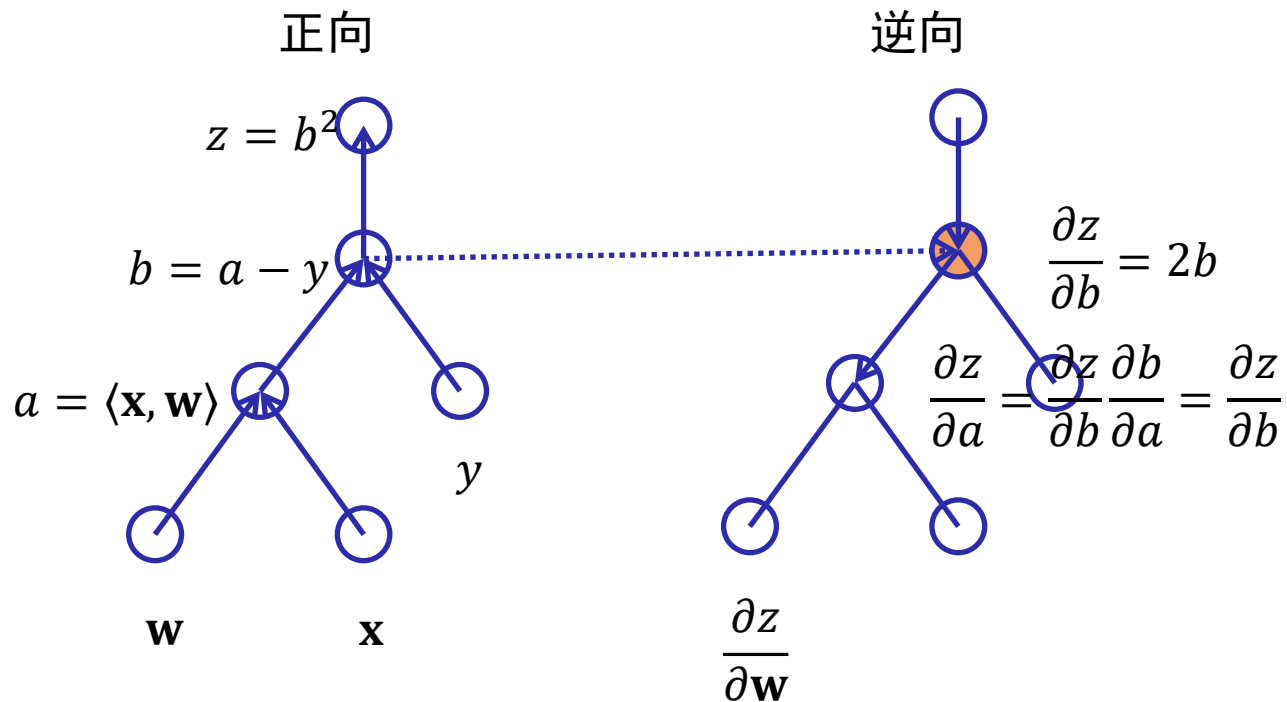
反向传播

► 假设 = $(\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$



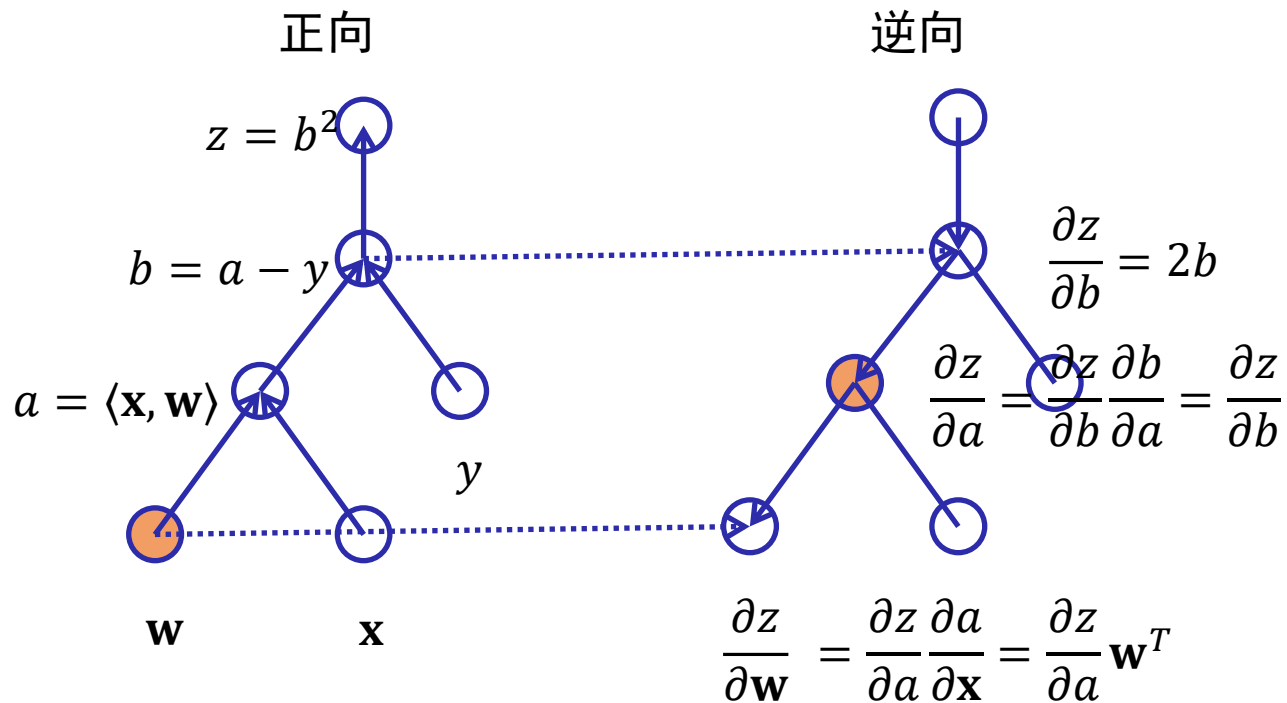
反向传播

► 假设 = $(\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$



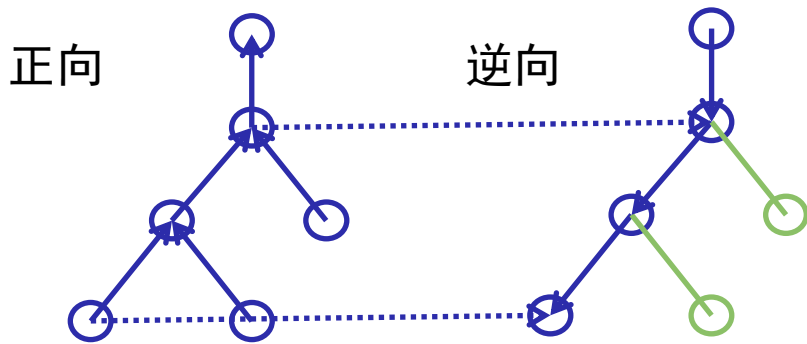
反向传播

► 假设 $= (\langle \mathbf{x}, \mathbf{w} \rangle - y)^2$



反向传播总结

- 创建一个计算图
- 正向：计算有向无环图，储存中间值
- 反向：逆向计算有向无环图
 - 减少不需要的图



复杂度

➤ $O(n)$, n 为计算次数

➤ 反向传播复杂度:

➤ 时间复杂度: $O(n)$, 计算所有导数, 基本上与正向复杂度一致

➤ 内存复杂度: $O(n)$, 需要储存所有正向计算的中间值

➤ 对比正向传播:

➤ 时间复杂度: $O(n)$, 计算 k 个变量的导数为 $O(n*k)$

➤ 内存复杂度: $O(1)$

[拓展] 再具体化 (re-materialization)

- 内存是逆向传播的瓶颈
 - 随着层数和批量大小线性增长
 - 有限 GPU 内存 (最多32GB)
- 用算力换内存
 - 只保存一部分中间计算值
 - 当需要时重新计算未保存中间值

再具体化 (re-materialization)

正向

逆向

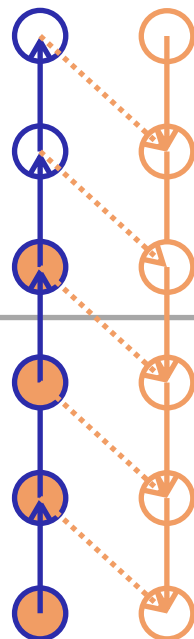
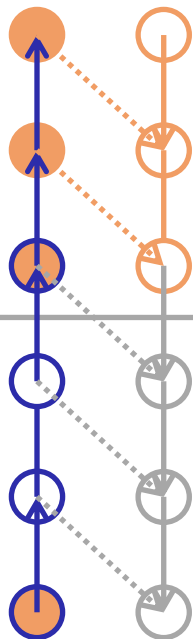
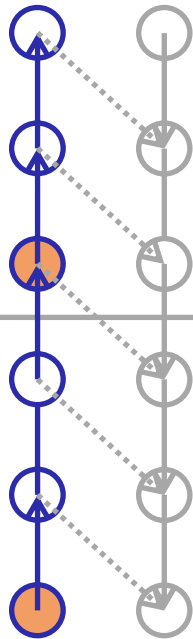
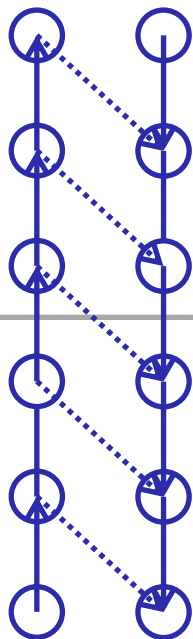
只保留每部分的
头部结果

重新计算第2部分
未保留的中间值

重新计算第1部分
未保留的中间值

第2部分

第1部分



复杂度

- 多一步正向传播
- 假设共有 m 部分, 则有 $O(m)$ 头部结果, 每个部分需内存 $O(n/m)$
 - 令 $m = \sqrt{n}$, 则内存复杂度为 $O(\sqrt{n})$
- 运用到深度学习网络
 - 只丢弃简单层, 如激活函数层, 常见 $<30\%$
 - 训练 10 倍大的网络, 或者 10 倍大的批量大小

总结

- 矩阵微积分
- 链式法则
- 自动微分法
- 反向传播