

20A\_Q2

April 19th /

## III Automatic Differentiation

$$l(x) = f(g(h(x)))$$



(a) reverse mode calculate  $\frac{dl}{dx}$  as follows,

at first forward and calc  $l(x)$

then backward from left to right

$$\frac{dl}{dx} = \underbrace{\frac{dl}{df} \frac{df}{dg} \frac{dg}{dh} \frac{dh}{dx}}$$

(b) forward mode evaluate each node and Jacobian as follows

$$\frac{dl}{dx} = \underbrace{\frac{dh}{dx} \frac{dg}{dh} \frac{df}{dg} \frac{dl}{df}}_{\rightarrow \text{from left to right}}$$

20A - Q2

2

(c) if  $m \gg n$  which mode?

$$x \in \mathbb{R}^n \quad l: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

forward mode is preferable

because its computational complexity

$$\text{is } O(n|E| + n|V|)$$

where  $|E|$  : num of edge

$|V|$  : num of vertex or cell.

Backward mode  $O(m|E| + m|V|)$



20A - Q2

3

## ② Exploding Gradients

(a)  $f^k(x) = \underbrace{(f_0 f_0 \dots f)}_{k\text{ times}}(x) = W^k x$

$$W^k x = \underbrace{(Q A Q^{-1}) \dots (Q A Q^{-1})}_{k\text{ times}} x$$

$$= Q A^k Q^{-1} x \in \mathbb{R}^n$$

$$= Q \begin{bmatrix} \lambda_1^k & & \\ & \ddots & \\ & & \lambda_n^k \end{bmatrix} Q^{-1} x \in \mathbb{R}^n$$

$k \rightarrow \infty$

if  $|\lambda_i| (1 \leq i \leq n) < 1$

$\rightarrow$  vanishing

$|\lambda_i| > 1$

$\rightarrow$  exploding

$\lambda_i = 1$

$\lambda_i = -1$

$\rightarrow$  converge

flip,  
not converge.

Do A - Q2

4

[2](b)

$$\frac{\partial f^k}{\partial x} = \frac{\partial f_k}{\partial f_{k-1}} \frac{\partial f_{k-1}}{\partial f_{k-2}} \cdots \frac{\partial f_1}{\partial x}$$



---

e.g.  $f(x) = 3x + 4$

$$\frac{\partial f^3}{\partial x} = \frac{\partial f(f(f(x)))}{\partial x}$$

$$\begin{aligned}f^3(x) &= f(f(f(x))) \\&= f(f(3x + 4)) \\&= f(9x + 12 + 4) = f(9x + 16) \\&= 27x + 16\end{aligned}$$

$$\frac{\partial f^3}{\partial x} = 27$$

---

20A - Q2

5

e.g.  $f(x) = 7x$

$$\begin{aligned}f^3(x) &= f(f(f(x))) \\&= f(f(7x)) \\&= f(49x) \\&= 7^3 x.\end{aligned}$$

$$\frac{\partial f^3}{\partial x} = 7^3$$



$$\begin{aligned}\frac{\partial f^k}{\partial x} &= \frac{\partial k}{\partial 0} = \frac{\partial k}{\partial (k-1)} \frac{\partial (k-1)}{\partial (k-2)} \dots \frac{\partial 1}{\partial 0} \\&= \prod_{i=1}^k w_i = w^k\end{aligned}$$

$$W = w_i \quad i \quad (1 \leq i \leq k)$$

20A-Q2 ( $f$  is not linear)

6

$\boxed{2} \text{ :: } A.$

$$\bar{x} + f u \quad u^T J = \lambda u^T$$
$$J \equiv \frac{\partial f}{\partial x}$$

$$\frac{\partial f}{\partial x} = J$$

$$(\bar{x} + f u)^T \frac{\partial f^k}{\partial x}(x) = (\bar{x} + f u)^T J^k$$

when we start from  $\bar{x}^T$ ,  $\bar{x}^T \frac{\partial f^k}{\partial x} = \bar{x}^T J^k$

here :  $f u^T J^k = f \lambda u^T J^{k-1}$   
 $= f \lambda^k u^T$

B. Why  $|\lambda_{\max}|$  plays crucial role of RNN

from here, since small perturbation

affects so much

if  $\lambda_{\max} > 1$



Accumulated gradient is affected  
by  $\lambda_{\max}$  in above way.

20A-Q2

7

③ Second-Order

$$\hat{f}_{x_0}(x) = f(x_0) + (x - x_0)^T g + \frac{1}{2} (x - x_0)^T H (x - x_0)$$

where  $g = \frac{\partial f}{\partial x}(x_0)$

$$H = \frac{\partial^2 f}{\partial x^2}(x_0)$$

(a)  $\hat{f}_{x_0}(x_0 - \varepsilon g) = f(x_0) - \varepsilon g^T g + \frac{\varepsilon^2}{2} g^T H g$

(b)  $\hat{f}_{x_0}(x_0 - \varepsilon g) = \hat{f}_{x_0}(\varepsilon) = \frac{1}{2} \varepsilon^2 g^T H g - \varepsilon g^T g + f(x_0)$

$$\varepsilon^* = \underset{\varepsilon}{\operatorname{argmin}} \hat{f}_{x_0}(\varepsilon)$$

↑  
and      is convex.

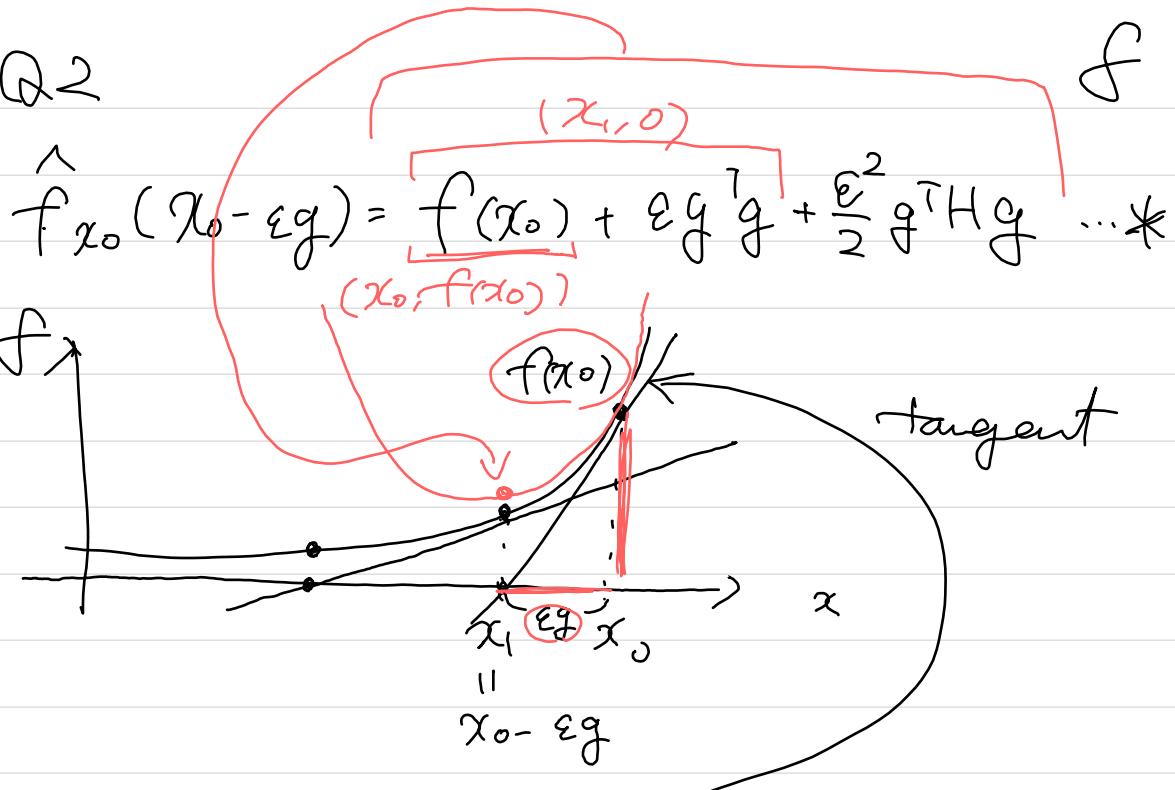
$$\text{So } \varepsilon^* = \left\{ \varepsilon \mid \frac{\partial}{\partial \varepsilon} \hat{f}_{x_0}(\varepsilon) = 0 \right\}$$

$$\Leftrightarrow \varepsilon g^T H g + g^T g = 0$$

$$\begin{aligned} \Leftrightarrow \varepsilon &= - (g^T H g)^{-1} g^T g \\ &= - H^{-1} \end{aligned}$$

20A - Q2

3) b



$$y = 0 = f(x_0) + f'(x_0)(x_1 - x_0)$$

$$\Leftrightarrow x_1 - x_0 = -\frac{f(x_0)}{f'(x_0)}$$

$$x_1 = x_0 + \frac{f(x_0)}{f'(x_0)}$$

$g$  is slope of tangent of  $f$  at  $x_0$

$H$  is curvature of  $f$  at  $x_0$

Answer

\* is convex approximation

where  $\gamma = x_0 - \varepsilon g$

# 20A - Q2

9

3(c)

$$\hat{f}_{x_0}(x) \equiv f(x_0) + (x - x_0)^T g + \frac{1}{2} (x - x_0)^T H (x - x_0)$$

$$\hat{f}_{x_0}(x_0 + \delta) = f(x_0) + f^T g + \frac{1}{2} f^T H f$$

$$\hat{f}_{x_0}(\delta) = \dots$$

$$f_{\text{opt}} = \underset{\delta}{\operatorname{argmin}} \hat{f}_{x_0}(\delta)$$

$$\frac{\partial \hat{f}_{x_0}(\delta)}{\partial \delta} = H\delta + g = 0$$

$$\Leftrightarrow \delta = -H^{-1}g$$

Newton update direction.

$$\hat{f}_{x_0}(x) = f(x_0) - \varepsilon \delta^T g + \frac{1}{2} \delta^T H \delta$$

$\downarrow$

$$= f(x_0) + \frac{\varepsilon^2}{2} \delta^T H \delta$$

if  $H \succeq 0$   $\varepsilon$  should be 0

//