

Automatic differentiation:

2-layer-NN:

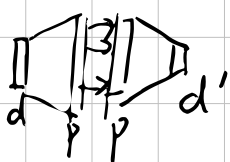
$$a \rightarrow \bigcup a \xrightarrow{g(a, \theta)} \sigma(\bigcup a) \rightarrow V \sigma(\bigcup a)$$

$g(a, \theta)$
“(U, V)”

$$\sigma(z) = (\sigma(z_i))_{i=1}^n$$

$$\sigma(s) = \frac{e^s}{1+e^s}$$

$$s(s) = \max(s, 0)$$



$$U = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \end{pmatrix}$$

$$V = v_1 v_2 \dots$$

$$g(a, \theta) = V \sigma(Ua) = \sum_{i=1}^p \underbrace{v_i}_{\text{"a neuron"}} \cdot \sigma(\underbrace{\langle u_i, a \rangle}_{\text{"a neuron"}})$$

$$\min_{\theta=U, V} \sum_{i=1}^n (g(a_i, \theta) - y_i)^2$$

$$\min_{U, V} \sum_{i=1}^n (\underbrace{V \sigma(Ua_i)}_{\text{easy } b_i} - y_i)^2$$

$b_i \approx f(U, V)$

Exercise: Compute $\nabla_U f$ $\nabla_V f$

$$Q^1: f: \mathbb{R}^d \rightarrow f(x) \in \mathbb{R}$$

Hyp: apply f in cost k

$$\nabla f: \mathbb{R} \rightarrow \mathbb{R}$$

$$\nabla f: \mathbb{R}^d \rightarrow \nabla f(x) \in \mathbb{R}^d$$

Q^2 : How many operations to compute $\nabla f(x)$

$$\nabla f(x) \approx \left(\frac{f(x+\epsilon e_1) - f(x)}{\epsilon}, \frac{f(x+\epsilon e_2) - f(x)}{\epsilon}, \dots \right) \partial_1 \dots \partial_d$$

$\hookrightarrow (d+1)k$ operations

Thm: cost of $\nabla f(x) \sim 3k$

$$+ \text{ algo } \underset{k}{\text{ Algo } (f)} \rightarrow \underset{3k}{\text{ Algo } (\nabla f)}$$

$$f(x) = f_k \circ f_{k-1} \circ \dots \circ f_0(x)$$

$$\nabla f(x) = \nabla f_k(x_k) \times \nabla f_{k-1}(x_{k-1}) \times \dots \times \nabla f_0(x_0=x)$$

$$x_0 = x$$

$$x_{k+1} = f_k(x_k) \quad x_k \in \mathbb{R}^{d_k}$$

$$\nabla f(x) = (A_k \times (A_{k-1} \times (A_{k-2} \dots))) \times A_0$$

Forward method: $A_k (A_{k-1} (A_{k-2} (A_{k-3} (A_0 \dots$

Backward $(\underbrace{A_k A_{k-1}}_{n^2}) A_{k-2} \dots A_0$

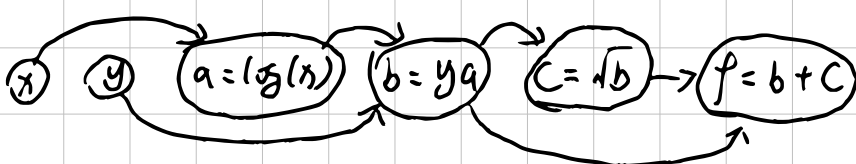
Ex. $d_k = 1 \quad d_{k-1} = d_{k-2} = \dots = d$ forward $n^2 + (k-1)n^3$
 $\hookrightarrow kn^2$

Computational graph

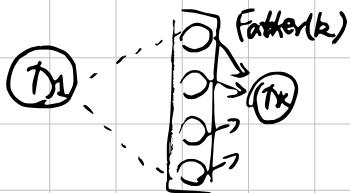
DAG: directed acyclic graph

Algo: $\gamma_k = f_k(\gamma_{k-2}, \gamma_{k-1}, \dots, \gamma_1)$ exam

$$f(x, y) = y \log(x) + \sqrt{y \log(x)}$$



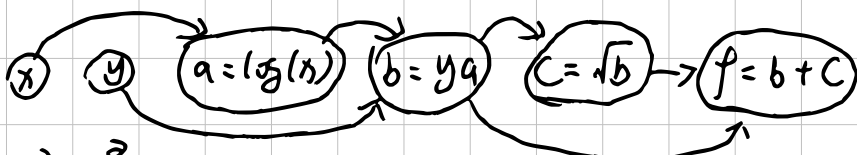
Forward = classical chain rule



$$\gamma_k = f_k((\gamma_i)_{i \in \text{Father}(k)})$$

$$\text{Thm: } \frac{\partial \gamma_k}{\partial \gamma_1} = \sum_{l \in \text{Father}(k)} \left[\frac{\partial \gamma_k}{\partial \gamma_l} \right] \times \frac{\partial \gamma_l}{\partial \gamma_1} \Rightarrow \frac{\partial f_k(\dots)}{\partial \gamma_1}$$

$$\text{Initia: } \frac{\partial \gamma_1}{\partial \gamma_1} = Id_{1 \times 1} \quad \frac{\partial \gamma_2}{\partial \gamma_1} = 0 \quad \dots$$



$$\text{FW: } \frac{\partial}{\partial x}$$

$$\frac{\partial x}{\partial x} = 1 \quad \frac{\partial y}{\partial x} = 0$$

$$\frac{\partial a}{\partial x} = \left[\frac{\partial a}{\partial x} \right] \times \frac{\partial x}{\partial x} = 1 \times \frac{\partial x}{\partial x}$$

$$\frac{\partial b}{\partial x} = \frac{\partial b}{\partial y} \frac{\partial y}{\partial x} + \frac{\partial b}{\partial a} \frac{\partial a}{\partial x} = a \frac{\partial y}{\partial x} + y \frac{\partial a}{\partial x}$$

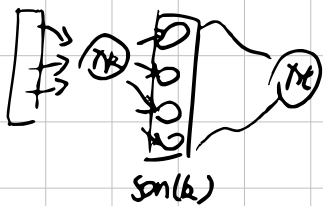
$$\frac{\partial c}{\partial x} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial x} = \frac{1}{2\sqrt{b}} \times \frac{\partial b}{\partial x}$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial b} \frac{\partial b}{\partial x} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial x} = 1 \frac{\partial b}{\partial x} + 1 \frac{\partial c}{\partial x}$$

$$\frac{\partial}{\partial y} \dots$$

$$\text{BWD } \frac{\partial \gamma_k}{\partial \gamma_k} \leftarrow \text{last} \quad \frac{\partial \gamma_k}{\partial \gamma_k} \leftarrow \text{current}$$

$$\text{FWD } \frac{\partial \gamma_k}{\partial \gamma_k} \leftarrow \text{current} \quad \frac{\partial \gamma_k}{\partial \gamma_k} \leftarrow \text{last}$$



$$\text{Thm: } \frac{\partial \gamma_k}{\partial \gamma_k} = \sum_{l \in \text{Son}(k)} \frac{\partial \gamma_l}{\partial \gamma_k} \left[\frac{\partial \gamma_l}{\partial \gamma_k} \right]$$

$$\frac{\partial \gamma_k}{\partial \gamma_k} = Id$$

In practice: $\gamma_k \in \mathbb{R}$ (last) f = \gamma_k

$$\frac{\partial f}{\partial \gamma_k} = \boxed{} = \left[\nabla_{\gamma_k} f \right]^T$$

$$\nabla_{\gamma_k} f = \sum_{l \in \text{Son}(k)} \left[\frac{\partial f_l}{\partial \gamma_k} \right]^T (\nabla_{\gamma_l} f)$$

$$\begin{matrix} \times & \boxed{} \\ \times & \boxed{}^T \end{matrix} \left. \vphantom{\begin{matrix} \times \\ \times \end{matrix}} \right\} \begin{matrix} \text{Jacobian} \\ \text{vector product} \end{matrix}$$

Summary:

Algo ($\gamma_1 - \gamma_b$)

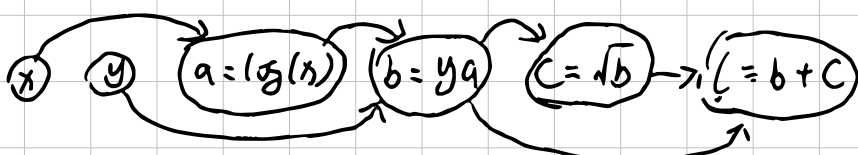
For $k = b+1 \dots t$

$$\left[\begin{array}{l} \gamma_k = f_k(\gamma_1 - \gamma_{k-1}) \\ \text{[Store } \gamma_k \text{]} \end{array} \right.$$

Backprop ($\gamma_1 - \gamma_b$) $\nabla_{\gamma_k} \ell = 1$

For $k = t-1, t-2, \dots, 1$

$$\nabla_{\gamma_k} \ell = \sum_{l \in \text{Son}(k)} \left[\frac{\partial f_l}{\partial \gamma_k} \right]^T (\nabla_{\gamma_l} \ell)$$



$$\frac{\partial \ell}{\partial \ell} = 1$$

$$\frac{\partial \ell}{\partial c} = \frac{\partial \ell}{\partial \ell} \times \left[\frac{\partial \ell}{\partial c} \right] = \frac{\partial \ell}{\partial \ell} \times 1$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial c} \left[\frac{\partial c}{\partial b} \right] + \frac{\partial \ell}{\partial (b+c)} \left[\frac{\partial \ell}{\partial b} \right] = \frac{\partial \ell}{\partial c} \times \frac{1}{2\sqrt{b}} + \frac{\partial \ell}{\partial \ell} \times 1$$

$$\frac{\partial \ell}{\partial a} = \frac{\partial \ell}{\partial b} \left[\frac{\partial b}{\partial a} \right] = \frac{\partial \ell}{\partial b} y$$

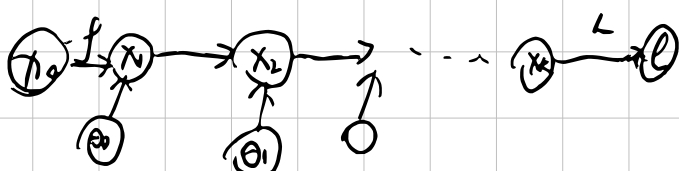
$$\frac{\partial \ell}{\partial y} = \frac{\partial \ell}{\partial b} \left[\frac{\partial b}{\partial y} \right] = \frac{\partial \ell}{\partial b} a$$

$$\frac{\partial \ell}{\partial x} = \frac{\partial \ell}{\partial a} \left[\frac{\partial a}{\partial x} \right] = \frac{\partial \ell}{\partial a} \frac{1}{x}$$

MLP: x_0 input (ex. image)

$$k=0 \dots k-1 \quad x_{k+1} = \sigma(\theta_k x_k)$$

$$f(x, \theta) = L(x_k)$$



$$f_1 = f_2 = \dots = f_k \quad (\text{layers})$$

$$x_{k+1} = f_{k+1}(x_k, \theta_k) = \sigma(\theta_k x_k)$$

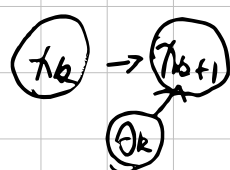
Back-prop $\nabla_{x_k} f$

$$\nabla_{\theta_k} f$$

$$\text{WIT: } \nabla_{x_k} f = \nabla L(x_k)$$

$$\text{Recursion: } \nabla_{x_k} f = \left[\frac{\partial f_k}{\partial x_k} \right]^T (\nabla_{x_{k+1}} f)$$

$$\nabla_{\theta_k} f = \left[\frac{\partial f_k}{\partial \theta_k} \right]^T (\nabla_{x_{k+1}} f)$$



$$f: \underbrace{(x, \theta)}_{\substack{\mathbb{R}^{d \times d} \\ \mathbb{R}^{d \times N}}} \mapsto \underbrace{\sigma(\theta x)}_{d \times N} \in \mathbb{R}^{d \times N} \quad (\text{setting})$$

$$\text{What is implementation } \frac{\partial f}{\partial x}(x, \theta)^T [z] = ??$$

$$\frac{\partial f}{\partial \theta}(x, \theta)^T [z] = ??$$

$$\left\{ \begin{array}{l} \frac{\partial f}{\partial x}(x, \theta): \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N} \text{ linear} \\ \frac{\partial f}{\partial x}(x, \theta)^T \in \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N} \end{array} \right.$$

$$\left\{ \begin{array}{l} \frac{\partial f}{\partial \theta}(x, \theta): \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times N} \\ \frac{\partial f}{\partial \theta}(x, \theta)^T \in \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times d} \end{array} \right.$$

$$f(x, \theta) = \sigma(\theta x)$$

$$f(x + \varepsilon T, \theta) = \sigma(\theta x + \varepsilon \theta T)$$

$$\sigma(w + \varepsilon \gamma) = \begin{pmatrix} \sigma(w_1 + \varepsilon \gamma_1) \\ \vdots \\ \sigma(w_d + \varepsilon \gamma_d) \end{pmatrix}$$

$$= \begin{pmatrix} \sigma(w_1) + \varepsilon \sigma'(w_1) \gamma_1 \\ \vdots \\ \sigma(w_d) + \varepsilon \sigma'(w_d) \gamma_d \end{pmatrix} + o(\varepsilon)$$

$$= \sigma(w) + \varepsilon \underbrace{\sigma'(w) \gamma}_{\text{diag}(\sigma'(w)) \times \gamma} + o(\varepsilon)$$

$$\text{Recap: } \partial \sigma(w) = \text{diag}(\sigma'(w))$$

$$\text{Ex. Relu}(s) = \frac{\sigma(s)}{s}$$

$$f(x, \theta) = \sigma(\theta x)$$

$$f(x + \varepsilon T, \theta) = \sigma(\theta x + \varepsilon \theta T)$$

$$= \sigma(\theta x) + \varepsilon \text{diag}(\sigma'(\theta x)) \times (\theta T) + o(\varepsilon)$$

\uparrow Jacob. of σ

$$\frac{\partial f}{\partial x}(x, \theta) [T] = \text{diag}(\sigma'(\theta x)) \times \theta T$$

$$\text{Remember: } \langle A [T], z \rangle = \langle T, A^T [z] \rangle$$

$$\left\langle \frac{\partial f}{\partial x}(x, \theta) [T], z \right\rangle = \langle \text{diag}(\sigma'(\theta x)) \times \theta T, z \rangle$$

$$= \langle T, \underbrace{\theta^T \text{diag}(\sigma'(\theta x))}_{\text{adjoint}} z \rangle$$

$$\text{Concl}^\circ 1: \frac{\partial f}{\partial x}(x, \theta)^T [z] = \theta^T \text{diag}(\sigma'(\theta x)) z$$

$$\text{on } \theta: f(x, \theta + \varepsilon z) = \sigma((\theta + \varepsilon z)x)$$

$$= \sigma(\theta x + \varepsilon z x) = \sigma(\theta x) + \varepsilon \text{diag}(\sigma'(\theta x)) \theta x + o(\varepsilon)$$

$$\frac{\partial f}{\partial \theta}(x, \theta) [z] = \text{diag}(\sigma'(\theta x)) \times z x$$

$$\left\langle \frac{\partial f}{\partial \theta}(x, \theta) [z], T \right\rangle = \langle \text{diag}(\dots) z x, T \rangle$$

$$= \langle z, \text{diag}(\dots)^T x^T \rangle$$

$$\text{Conclu}^\circ 2: \frac{\partial f}{\partial \theta}(x, \theta)^T [z] = \underbrace{\text{diag}(\sigma'(\theta x))}_{\mathbb{R}^{d \times N}} \underbrace{z x^T}_{\substack{\mathbb{R}^{d \times N} \times \mathbb{R}^{N \times d} \\ \mathbb{R}^{d \times d}}}$$

$$\theta_k = (U_k, V_k)$$

$$x_{k+1} = x_k + U_k \sigma(V_k x_k)$$

$$\text{Residual } x_k \xrightarrow{U} \xrightarrow{\sigma} \xrightarrow{V} z_k \xrightarrow{\mathcal{D}} x_{k+1}$$