

Proposal Due 4/8 11:59 PM

→ Literature review: no need to re implement

Neural Network Optimization

W

Machine Learning Problems

- **Given data:**

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- **Learning a model's parameters:** $\underset{w}{\underbrace{\sum_{i=1}^n \ell_i(w)}}$

Logistic Loss: $\ell_i(w) = \log(1 + \exp(-y_i x_i^T w))$

Squared error Loss: $\ell_i(w) = (y_i - x_i^T w)^2$

Machine Learning Problems

- Given data:

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- Learning a model's parameters:

$$\sum_{i=1}^n \ell_i(w)$$

Logistic Loss: $\ell_i(w) = \log(1 + \exp(-y_i x_i^T w))$

Squared error Loss: $\ell_i(w) = (y_i - x_i^T w)^2$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \left(\frac{1}{n} \sum_{i=1}^n \ell_i(w) \right) \Big|_{w=w_t}$$

η : step size
learning rate

w_0 : random init
Learn init, training init

Gradient Descent

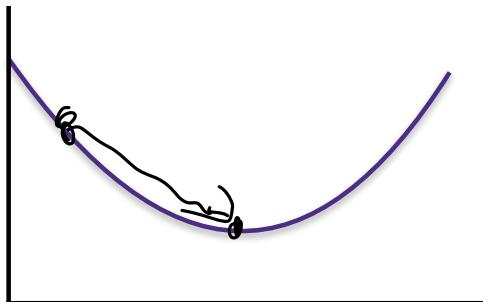
Initialize: $w_0 = 0$

for $t = 1, 2, \dots$

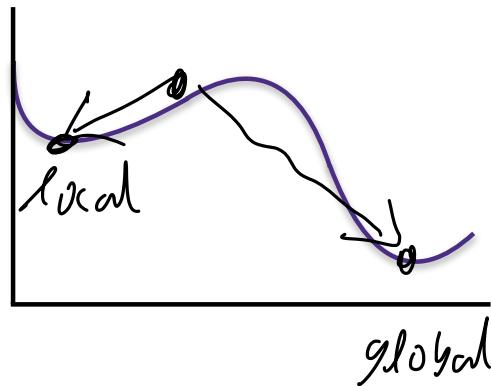
$$w_{t+1} = w_t - \eta \nabla f(w_t)$$



Convex Function



Non-convex Function



Sub-Gradient Descent

$\mathcal{L}_\ell L \cup$

Initialize: $w_0 = 0$

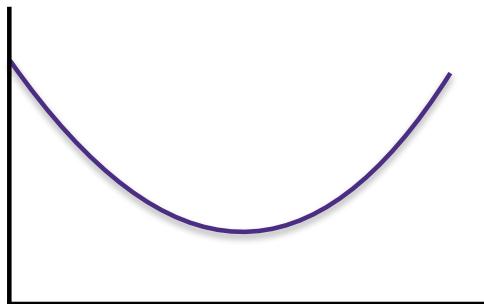
for $t = 1, 2, \dots$

Find any g_t such that $f(y) \geq f(w_t) + g_t^\top (y - w_t)$

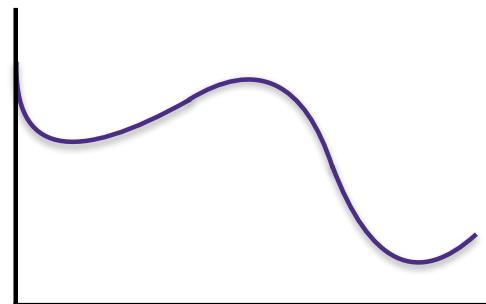
$$w_{t+1} = w_t - \eta g_t$$

g is a subgradient at x if $f(y) \geq f(x) + g^T(y - x)$

Convex Function



Non-convex Function



Machine Learning Problems

- **Given data:**

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- **Learning a model's parameters:**

$$\sum_{i=1}^n \ell_i(w)$$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \left(\frac{1}{n} \sum_{i=1}^n \ell_i(w) \right) \Big|_{w=w_t}$$

Stochastic Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w) \Big|_{w=w_t}$$

I_t drawn uniform at random from $\{1, \dots, n\}$

Mini-batch SGD

B: batch size

Instead of one iterate, average B stochastic gradient together

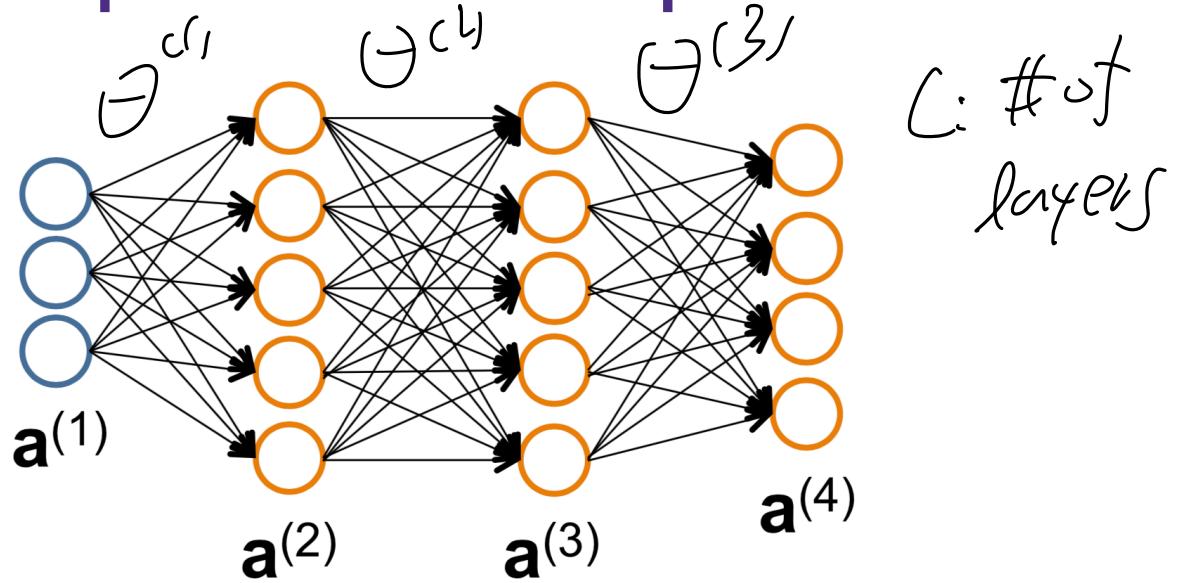
Advantages:

- de-noises gradient
- Matrix computations
- Parallelization

$$\frac{1}{B} \sum_{j=1}^B \nabla \ell(j) \quad (\text{w} \leftarrow \text{w})$$

Gradient Computation on a Graph

$$\frac{\partial L}{\partial \Theta^{(1)}}$$



Naive computation: node by node

$$\mathcal{O}(L^2)$$

A brief history

- **Back propagation:** the workhorse for training neural networks. An algorithm that for a network with V nodes and E edges calculates that gradient in **linear time** $O(V+E)$. $\mathcal{O}(V^2)$
- The name was introduced by Rumelhart, Hinton, Williams '86. Same idea was rediscovered multiple times. Also mentioned in Werbos' thesis '74 in the context of neural networks.
- **Control theory:** Kelly '60, Bryson '61 [**dynamic programming**]
- **Theoretical computer science:** Baur-Strassen lemma '83 [**algebraic circuits**]

$$a^{(1)} = x$$

$$\underline{z}^{(2)} = \underline{\Theta}^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

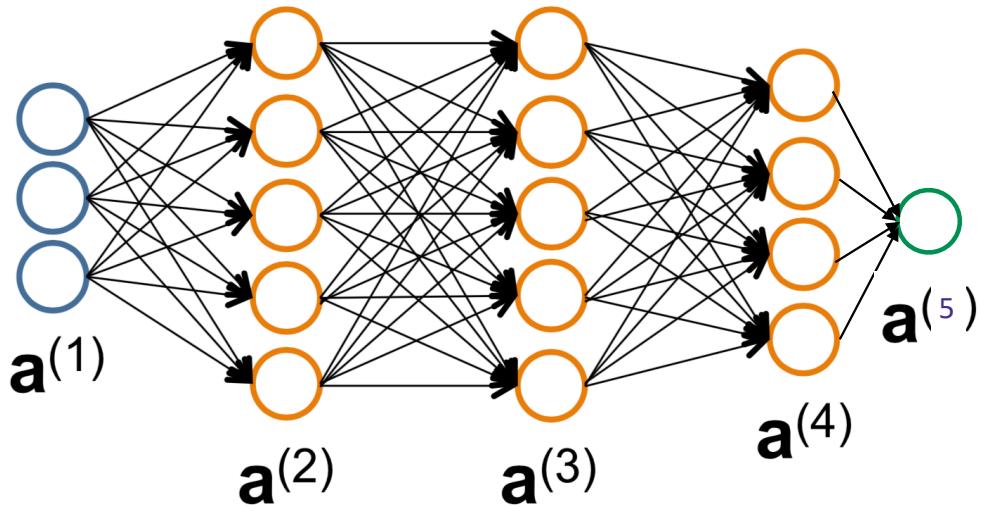
⋮

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\widehat{y} = g(\underline{\Theta}^{(L)} a^{(L)})$$



$$L(y, \widehat{y}) = y \log(\widehat{y}) + (1 - y) \log(1 - \widehat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Gradient Descent: $\Theta^{(l)} \leftarrow \Theta^{(l)} - \eta \nabla_{\Theta^{(l)}} L(y, \widehat{y}) \quad \forall l$

Forward Propagation

$$a^{(1)} = \underline{x}$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$\underline{z^{(l+1)}} = \Theta^{(l)} a^{(l)}$$

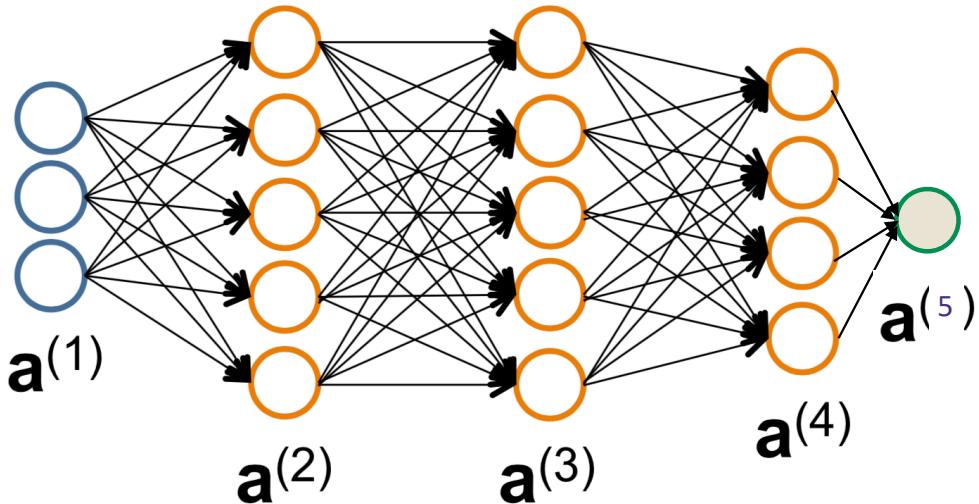
$$\underline{a^{(l+1)}} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

L : # of layers
ignore bias

g : activation function
 $z^{(l)}$: pre-activation



$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Backprop

$$a^{(1)} = x \in \mathbb{R}^d$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

$\Theta^{(1)}, \dots, \Theta^{(L)}$: parameters to train
 $\Theta^{(1)} \in \mathbb{R}^{m \times d}$, $\Theta^{(2)}, \dots, \Theta^{(L-1)} \in \mathbb{R}^{m \times m}$
 m : width $\Theta^{(L)} \in \mathbb{R}^m$

Train by Stochastic Gradient Descent:

$$\Theta_{i,j}^{(l)} \in \mathbb{R} \quad \Theta_{i,j}^{(l)} \leftarrow \Theta_{i,j}^{(l)} - \eta \frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}}$$

$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Backprop

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$\underline{z^{(l+1)}} = \underline{\Theta^{(l)} a^{(l)}}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

Chain Rule

$$z_i^{(l+1)} = \sum_{j=1}^m \theta_{i,j}^{(l+1)} \cdot a_j^{(l)}$$

$$\frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial \Theta_{i,j}^{(l)}} =: \delta_i^{(l+1)} \cdot a_j^{(l)}$$

maps $a_j^{(l+1)} \rightarrow z_i^{(l+1)}$

Train by Stochastic Gradient Descent:

$$\Theta_{i,j}^{(l)} \leftarrow \Theta_{i,j}^{(l)} - \eta \frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}}$$

$$\boxed{z_i^{(l+1)}} = \boxed{\Theta^{(l+1)}} \boxed{a^{(l)}}$$

$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Backprop

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

$$\frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial \Theta_{i,j}^{(l)}} =: \delta_i^{(l+1)} \cdot a_j^{(l)}$$

(Hierin versteht man die Rückwärtsausbreitung der Fehlerfunktion)

$$\delta_i^{(l)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l)}} = \sum_k \frac{\partial L(y, \hat{y})}{\partial z_k^{(l+1)}} \cdot \frac{\partial z_k^{(l+1)}}{\partial z_i^{(l)}}$$

$\delta_k^{(l+1)}$ ist die Ausbreitung des Fehlers von der Schicht $k+1$ zur Schicht k .
 $\Theta_{kj}^{(l)}$ ist die Gewichtung zwischen den Neuronen j in Schicht l und k in Schicht $k+1$.

$$z_R^{(l+1)} = \sum_{k=1}^m \Theta_{kj}^{(l)} \cdot \delta_k^{(l+1)} \cdot g'(z_k^{(l)})$$

$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Backprop

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

$$\frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial \Theta_{i,j}^{(l)}} =: \delta_i^{(l+1)} \cdot a_j^{(l)}$$

$$\begin{aligned}\delta_i^{(l)} &= \frac{\partial L(y, \hat{y})}{\partial z_i^{(l)}} = \sum_k \frac{\partial L(y, \hat{y})}{\partial z_k^{(l+1)}} \cdot \frac{\partial z_k^{(l+1)}}{\partial z_i^{(l)}} \\ &= \sum_k \delta_k^{(l+1)} \cdot \Theta_{k,i}^{(l)} g'(z_i^{(l)}) \\ &= a_i^{(l)}(1 - a_i^{(l)}) \sum_k \delta_k^{(l+1)} \cdot \Theta_{k,i}^{(l)}\end{aligned}$$

$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Backprop

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

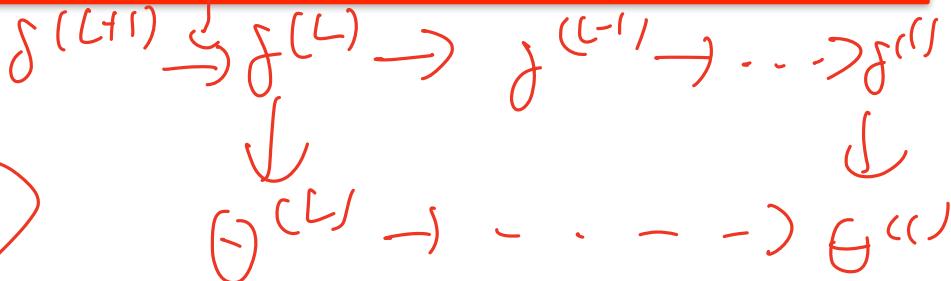
⋮

$$\hat{y} = a^{(L+1)}$$

Recursion / Dynamic Programming

$$\frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial \Theta_{i,j}^{(l)}} =: \delta_i^{(l+1)} \cdot a_j^{(l)}$$

$$\delta_i^{(l)} = a_i^{(l)}(1 - a_i^{(l)}) \sum_k \delta_k^{(l+1)} \cdot \Theta_{k,i}^{(l)}$$



$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Backprop

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

$$\frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial \Theta_{i,j}^{(l)}} =: \delta_i^{(l+1)} \cdot a_j^{(l)}$$

$$\delta_i^{(l)} = a_i^{(l)}(1 - a_i^{(l)}) \sum_k \delta_k^{(l+1)} \cdot \Theta_{k,i}^{(l)}$$

$$\begin{aligned}\delta_i^{(L+1)} &= \frac{\partial L(y, \hat{y})}{\partial z_i^{(L+1)}} = \frac{\partial}{\partial z_i^{(L+1)}} [y \log(g(z^{(L+1)})) + (1-y) \log(1-g(z^{(L+1)}))] \\ &= \frac{y}{g(z^{(L+1)})} g'(z^{(L+1)}) - \frac{1-y}{1-g(z^{(L+1)})} g'(z^{(L+1)}) \\ &= y - g(z^{(L+1)}) = y - a^{(L+1)}\end{aligned}$$

$$L(y, \hat{y}) = y \log(\hat{y}) + (1-y) \log(1-\hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Backprop

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

⋮

$$a^{(l)} = g(z^{(l)})$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$

$$\frac{\partial L(y, \hat{y})}{\partial \Theta_{i,j}^{(l)}} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial \Theta_{i,j}^{(l)}} =: \delta_i^{(l+1)} \cdot a_j^{(l)}$$

$$\delta_i^{(l)} = a_i^{(l)}(1 - a_i^{(l)}) \sum_k \delta_k^{(l+1)} \cdot \Theta_{k,i}^{(l)}$$

$$\delta^{(L+1)} = y - a^{(L+1)}$$

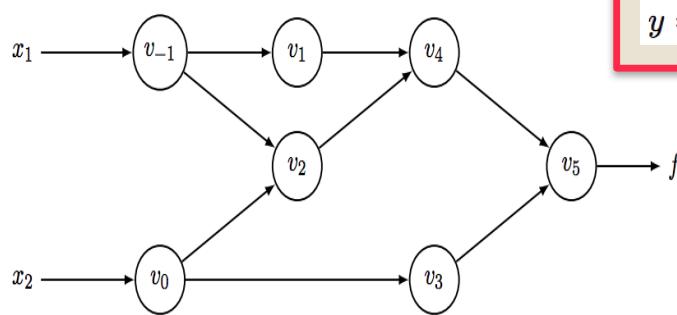
Recursive Algorithm!

$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \delta_i^{(l+1)} = \frac{\partial L(y, \hat{y})}{\partial z_i^{(l+1)}}$$

Auto-differentiation

Backprop for this simple network architecture is a special case of *reverse-mode auto-differentiation*:



$$y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$$

Forward Primal Trace	
$v_{-1} = x_1$	= 2
$v_0 = x_2$	= 5
$v_1 = \ln v_{-1}$	= $\ln 2$
$v_2 = v_{-1} \times v_0$	= 2×5
$v_3 = \sin v_0$	= $\sin 5$
$v_4 = v_1 + v_2$	= $0.693 + 10$
$v_5 = v_4 - v_3$	= $10.693 + 0.959$
$y = v_5$	= 11.652

Reverse Adjoint (Derivative) Trace	
$\bar{x}_1 = \bar{v}_{-1}$	= 5.5
$\bar{x}_2 = \bar{v}_0$	= 1.716
$\bar{v}_{-1} = \bar{v}_{-1} + \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}}$	= $\bar{v}_{-1} + \bar{v}_1/v_{-1} = 5.5$
$\bar{v}_0 = \bar{v}_0 + \bar{v}_2 \frac{\partial v_2}{\partial v_0}$	= $\bar{v}_0 + \bar{v}_2 \times v_{-1} = 1.716$
$\bar{v}_{-1} = \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}}$	= $\bar{v}_2 \times v_0 = 5$
$\bar{v}_0 = \bar{v}_3 \frac{\partial v_3}{\partial v_0}$	= $\bar{v}_3 \times \cos v_0 = -0.284$
$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2}$	= $\bar{v}_4 \times 1 = 1$
$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1}$	= $\bar{v}_4 \times 1 = 1$
$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$	= $\bar{v}_5 \times (-1) = -1$
$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$	= $\bar{v}_5 \times 1 = 1$
$\bar{v}_5 = \bar{y}$	= 1

Auto-differentiation

- Given a function, computes its partial derivatives
- Compute all of the partial derivatives of a function with (nearly) same computation runtime [Griewank '89, Baur and Strassen '83]
- Backbone of (applied) machine learning: Pytorch, Tensorflow, ...

Example of Computation Graph

$$f(w_1, w_2) = \left(\sin\left(\frac{2\pi w_1}{w_2}\right) + \underbrace{\frac{3w_1}{w_2} - \exp(2w_2)}_{z_4} \right) \cdot \left(\frac{3w_1}{w_2} - \exp(2w_2) \right)$$

Input: $z_0 = (w_1, w_2)$

1. $z_1 = w_1 \cdot w_2^{-1}$
2. $z_2 = \sin(2\pi z_1)$
3. $z_3 = \exp(2w_2)$
4. $z_4 = 3z_1 - z_3$
5. $z_5 = z_2 + z_4$
6. $z_6 = z_5 \cdot z_4$

Return z_6

The computation graph illustrates the flow of data from the input z_0 through various intermediate nodes to the final output z_6 . The nodes are represented by circles containing values or expressions. The edges represent the operations performed on the inputs to produce the outputs.

- Input:** $z_0 = (w_1, w_2)$
- Intermediate Nodes:**
 - $w_1 = [z_0]_1 \rightarrow z_1$
 - $w_2 = [z_0]_2 \rightarrow z_3$
 - $z_1 \rightarrow z_2$
 - $z_3 \rightarrow z_4$
 - $z_1 \rightarrow z_5$
 - $z_4 \rightarrow z_5$
 - $z_5 \rightarrow z_6$
- Final Output:** z_6

Computation Model

- Given access to a set of differentiable real functions $h \in \mathcal{H}$
- Use functions in \mathcal{H} to create intermediate variables.
- Evaluation trace:
 - All intermediate variables will be scalars; each corresponds to a node.
 - Input $z_0 = w \in \mathbb{R}^d$. $[z_0]_1 = w_1, [z_0]_2 = w_2, \dots, [z_0]_d = w_d$
 - Step 1: $z_1 = h_1$ (a subset of variables in w)
$$z_1 = h_1(w_1, w_2, w_3, w_4, w_5, w_6)$$
 - Step t: $z_t = h_t$ (a subset of variables in z_1, \dots, z_{t-1}, w)
 - ...
 - Step T: $z_T = h_T$ (a subset of variables in z_1, \dots, z_{T-1}, w)
 - **Return:** z_T
$$(h_1, \dots, h_T \in \mathcal{H})$$

Computation Model

- Every $h \in \mathcal{H}$ is one of the following:
 - Type 1: An affine transformation of the inputs

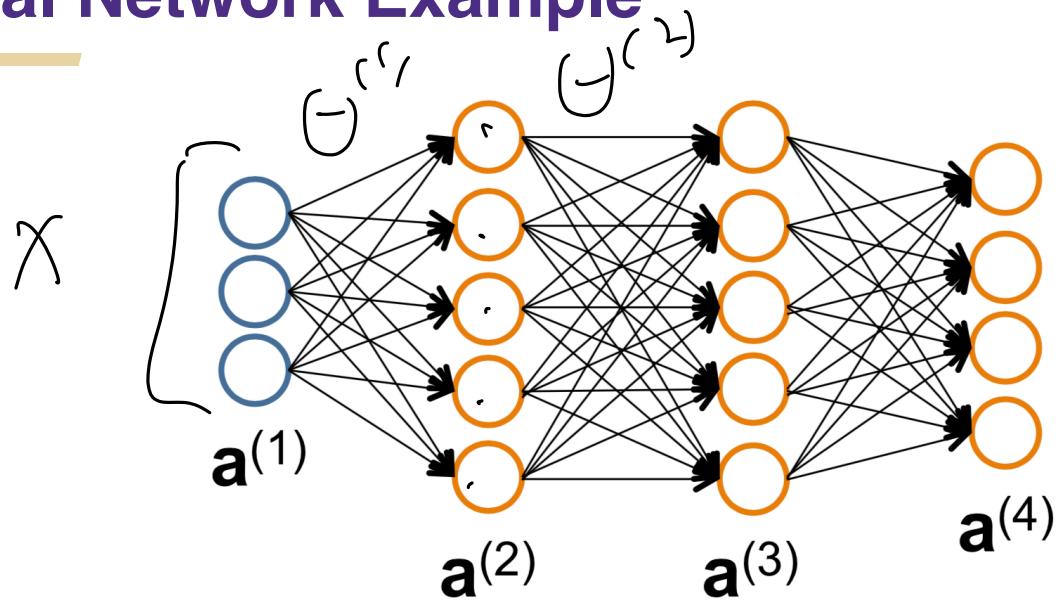
$$0.9 \cdot z_1 - z_3$$

- Type 2: A product of variables, to some power

$$w_1 \cdot w_2^{-1} \quad \text{or} \quad z_7 = z_1^4 z_5^5 z_5^{-4}$$

- Type 3: A fixed set of one dimensional differentiable functions: $\sin(\cdot)$, $\cos(\cdot)$, $\exp(\cdot)$, $\log(\cdot)$, ...
 - We assume we can easily compute the derivatives for each of these functions.
 - Type 3 can be approximated by Type 1 and Type 2, using polynomials.

Neural Network Example

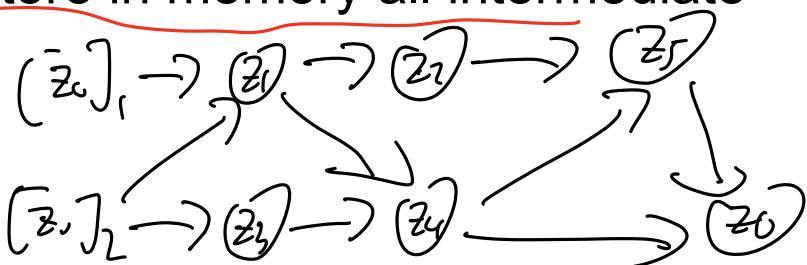


Reverse Mode of Automatic Differentiation

V : # of nodes E : # of edges

Goal: Compute partial derivatives of $f(w)$, i.e., df/dw .

- Step 1: computer $f(w)$ and store in memory all intermediate variables z_1, \dots, z_T



- Step 2: Initialize: $\frac{dz_T}{dz_T} = 1$.

- Step 3: For $t = T, T-1, \dots, 0$

$$\frac{dz_T}{dz_t} = \sum_{c \text{ is a child of } t} \frac{dz_T}{dz_c} \cdot \frac{\partial z_c}{\partial z_t}$$

(Child: a node z_t directly points to)

$$(1) \frac{dz_6}{dz_6} = 1$$

$$(2) \frac{dz_6}{dz_5} = \frac{dz_6}{dz_6} \cdot \frac{dz_5}{dz_5}$$

$$(3) \frac{dz_6}{dz_4} = \frac{dz_6}{dz_6} \cdot \frac{dz_5}{dz_4} + \frac{dz_6}{dz_5} \cdot \frac{dz_4}{dz_4}$$

$\mathcal{O}(V+E)$

- Step 4: Return $\frac{dz_T}{dz_0} = \frac{df}{dw}$

Time Complexity

Theorem (Baur and Strassen '83, Griewak '89): Assume every h is specified as in our computational model. For $h(\cdot)$ of type 3, assume we can compute the derivative $h'(z)$ in time as the same order of computing $h(z)$. Let T denote the time to compute $f(w)$. Then the reverse mode computes df/dw in time $O(T)$.

Pf: (1) Time: Count edges

(2) Correctness: $\frac{d z_c}{d z_\ell}$

$$\left\{ \begin{array}{l} \text{type 1: affine} \rightarrow \text{coefficient} \\ \text{type 2: product: } \frac{d z_r}{d z_\ell} = \frac{z_c - \alpha}{z_\ell} + \frac{d z_s}{d z_\ell} = \frac{z_s^2}{z_\ell} \\ \text{type 3: } z_c = h(z_\ell), \quad h'(z_\ell) \end{array} \right.$$

$z_5 = z_1 \cdot z_4$
 $\frac{d z_5}{d z_4} = \frac{z_5^2}{z_4}$

Time Complexity
