

15-780 – Graduate Artificial Intelligence: Neural networks

Aditi Raghunathan

Recap of supervised ML

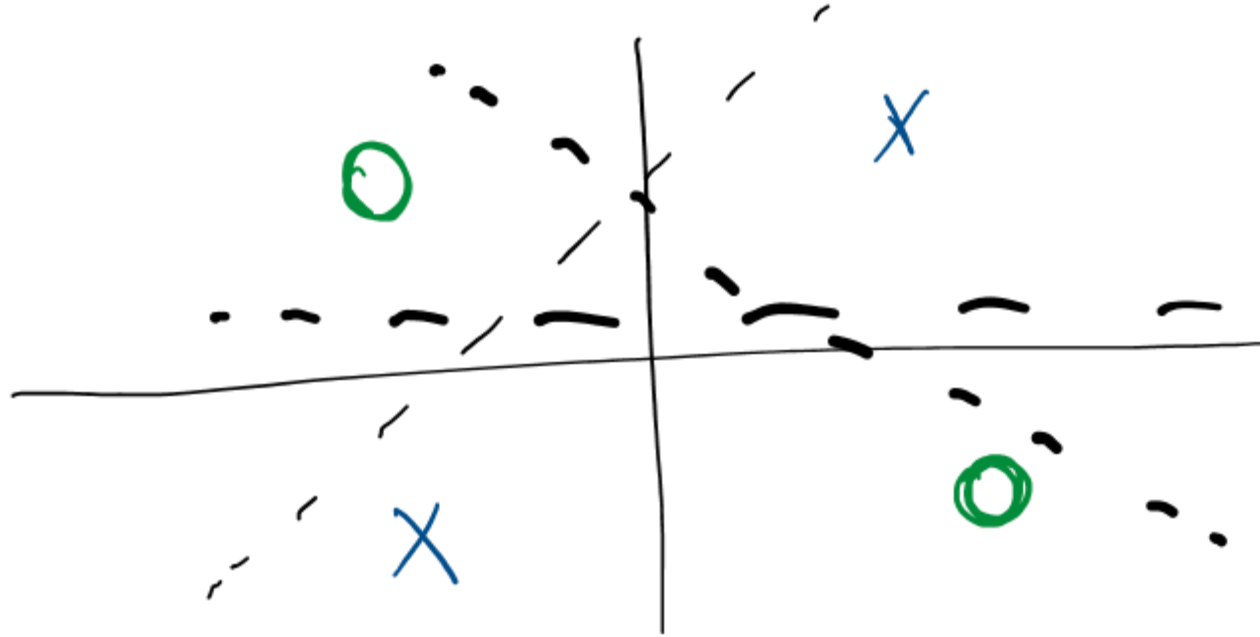
- Hypothesis function $h_{\theta}(x) = \theta^T x \Rightarrow \text{MLP}$
- Loss functions \rightarrow cross entropy for classification
- Optimization \rightarrow GD \rightarrow SGD \rightarrow variants like Adam

Linear classifiers

$$h_{\theta}(x) = \theta^T \underbrace{x}_{\text{diff features of } x}$$

$$= \theta^T \underbrace{\phi(x)}_{\text{carefully crafted features}}$$

Linear classifiers: XOR? → Minsky



Linear classifier -> "feature" learning

$$h_{\theta}(x) = \theta^T \underbrace{x}_{\text{new featurization}}$$

$$= \theta^T \underbrace{\phi(x)}_{\text{"hand-crafted"}}$$

→ Vision : SIFT, HOG features

→ kernels

Goal: learn features as well!

First attempt at feature learning

Goal: learn features

original input in \mathbb{R}^n

$$\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d \quad \text{linear}$$

$$h_{\theta}(x) = \theta^T \phi(x)$$

$$\theta \in \mathbb{R}^{d \times k}$$

$$\phi(x) = \theta_2^T x$$

$$h_{\theta}(x) = \theta_1^T \theta_2^T x$$

features are linear

$$h_{\theta}(x) = \tilde{\theta}^T x$$

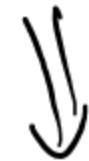
Second attempt

takeaway: need non-linearity

$$\phi(x) = \sigma(\theta_2^T x)$$

$$h_\theta(x) = \theta_1^T \sigma(\theta_2^T x)$$

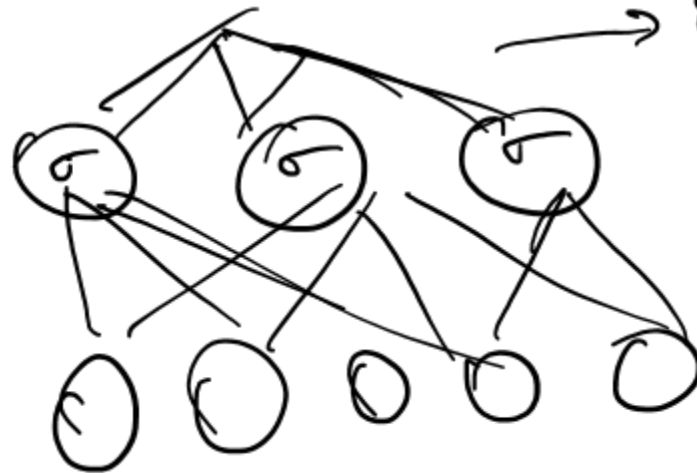
σ : non-linearity



activation fn

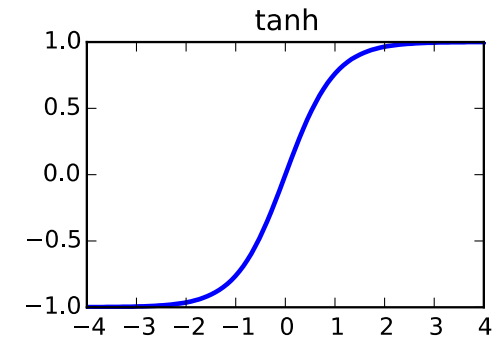
$\rightarrow \theta_1$ layer

$\theta_2^T x$ layer

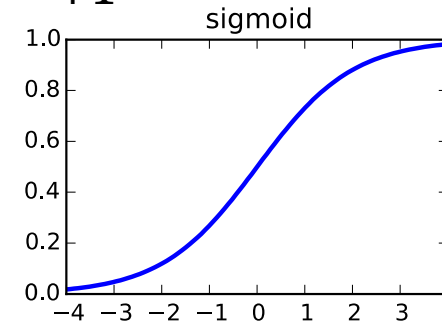


Activation functions

Hyperbolic tangent: $f(x) = \tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$

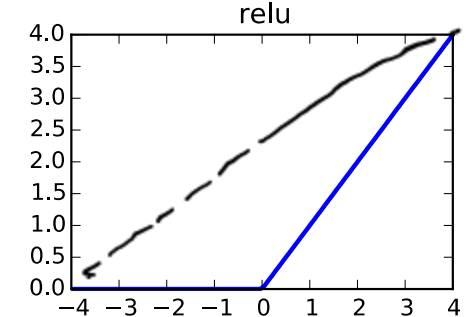


Sigmoid: $f(x) = \sigma(x) = \frac{1}{1+e^{-x}}$



Rectified linear unit (ReLU): $f(x) = \max\{x, 0\}$

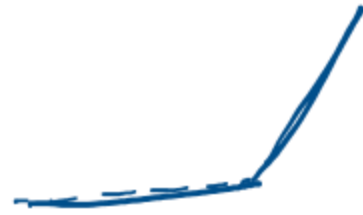
↓
off when $x < 0$
 x when $x > 0$



XOR example

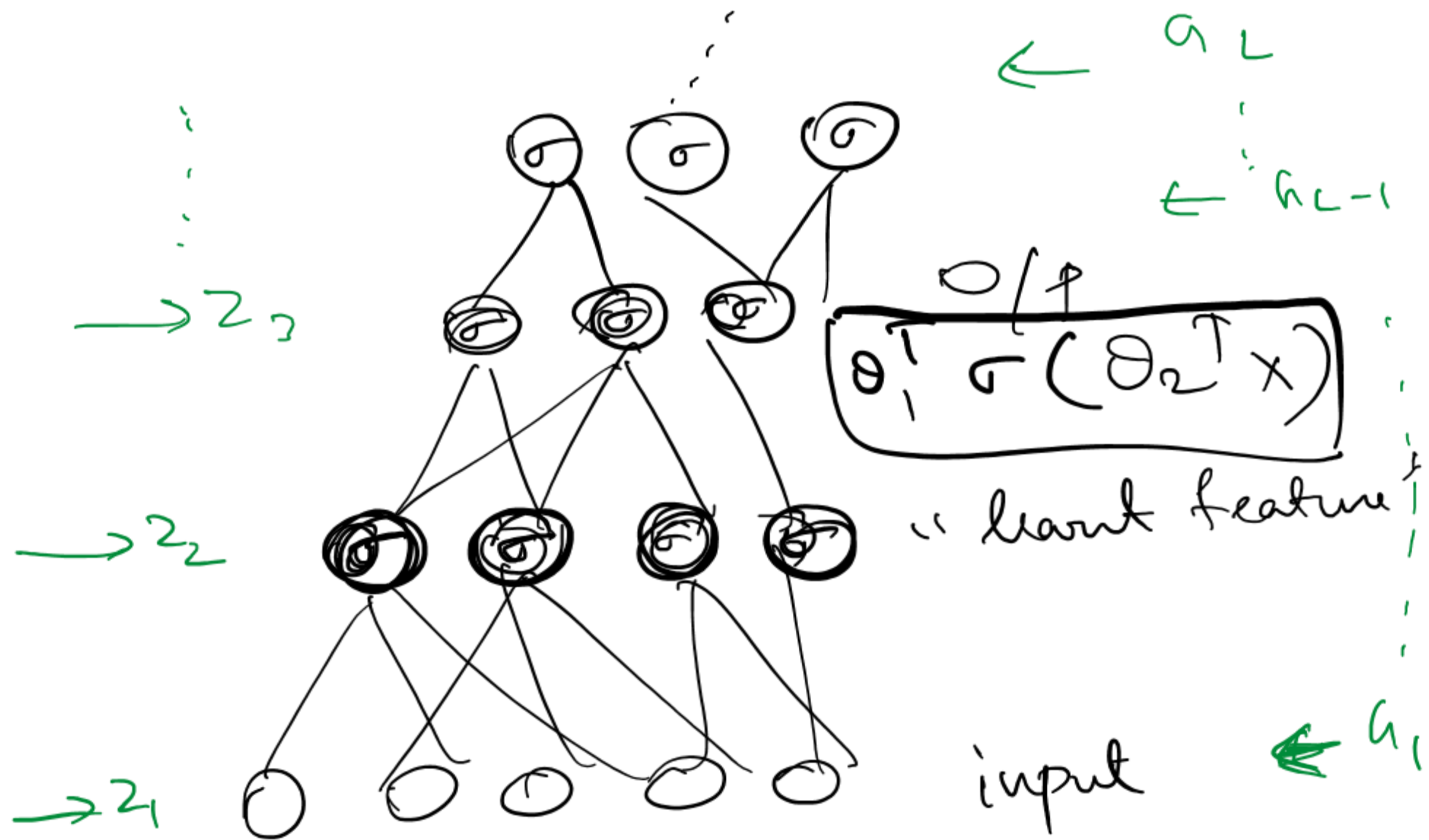
$$h_{\sigma}(x) = \theta_1^T \sigma(\theta_2^T x)$$

σ : ReLU



XOR: $\text{ReLU}(x_1 - x_2) + \text{ReLU}(x_2 - x_1)$

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0



"Deep" neural networks

- Multi-layer perceptron repeat the "hidden computation"

$$h_{\theta}(x) = w_L^T \left(\sigma(w_{L-1}^T \sigma(\dots \sigma(w_2^T \sigma(w_1^T x)) \dots) \right)$$

$$h_{\theta}(x) = w_2^T [\sigma(w_1^T x)]$$

parameters " θ ":
 $\{w_1, w_2 \dots w_L\}$

→ Fully connected network, feed-forward network

Universal function approximation

Theorem (1D case): Given any smooth function $f: \mathbb{R} \rightarrow \mathbb{R}$, closed region $\mathcal{D} \subset \mathbb{R}$, and $\epsilon > 0$, we can construct a one-hidden-layer neural network \hat{f} such that

$$\max_{x \in \mathcal{D}} |f(x) - \hat{f}(x)| \leq \epsilon$$

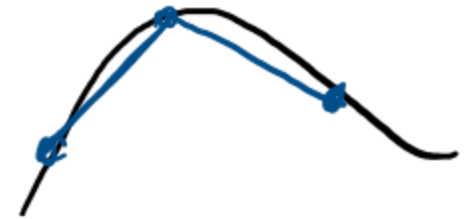
Proof: Select some dense sampling of points $(x^{(i)}, f(x^{(i)}))$ over \mathcal{D} . Create a neural network that passes exactly through these points (see below). Because the neural network function is piecewise linear, and the function f is smooth, by choosing the $x^{(i)}$ close enough together, we can approximate the function arbitrarily closely.

$$h(x) = \text{ReLU}(x - a)$$

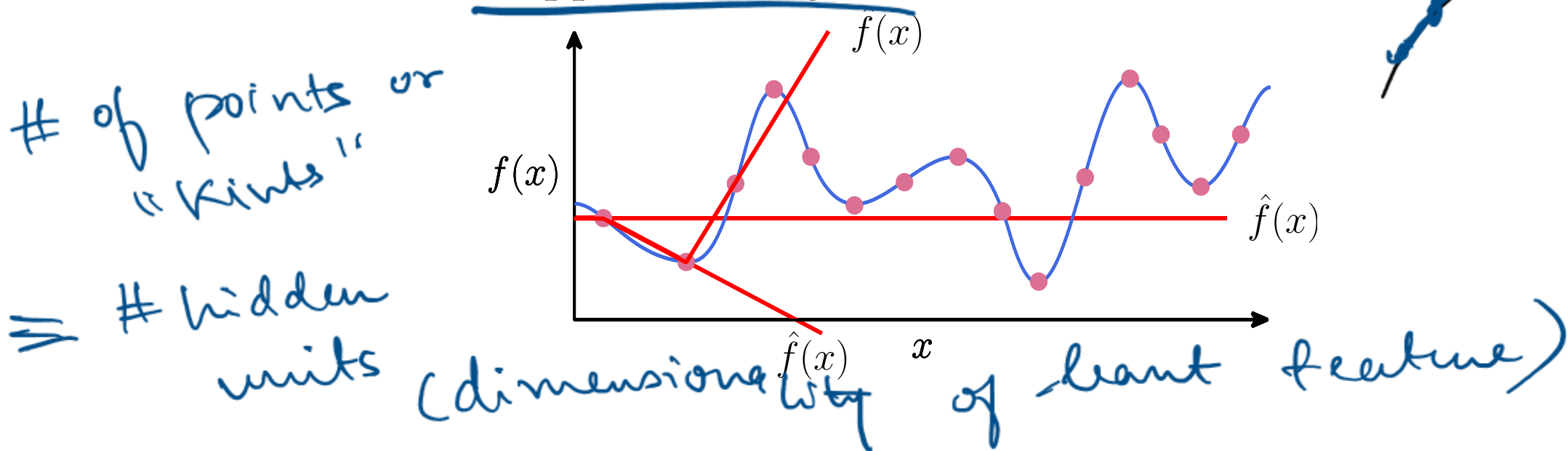
Universal function approximation

Assume one-hidden-layer ReLU network:

$$\hat{f}(x) = \sum_{i=1}^d \pm \max\{0, w_i x + b_i\}$$



Visual construction of approximating function.



Backpropagation

For SGD or variants, we need to compute **gradients of the loss** with respect to weights (parameters)

The gradient (recap)

A key concept in solving optimization problems is the notation of the gradient of a function (multi-variate analogue of derivative)

Derivative: $f'(x) = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$

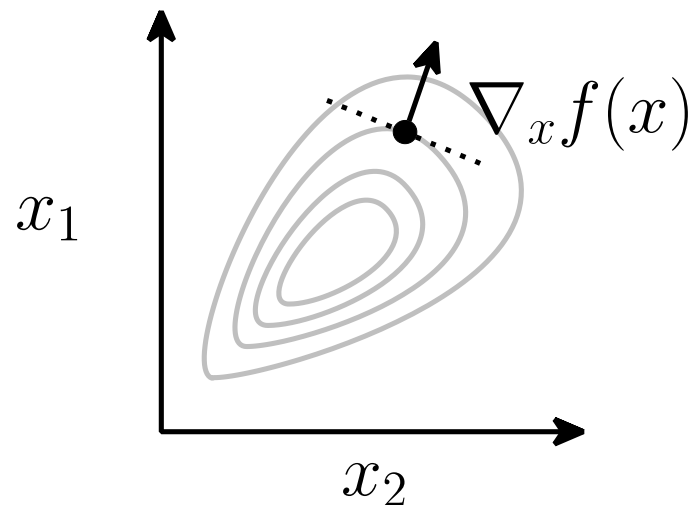
Partial derivative: A partial derivative of a function of several variables is derivative with respect to one of those variables with rest constant

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x + h\mathbf{e}_i) - f(x)}{h}$$

The gradient (recap)

For $f: \mathbb{R}^n \rightarrow \mathbb{R}$, gradient is defined as vector of partial derivatives

$$\nabla_x f(x) \in \mathbb{R}^n = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$



Points in “steepest direction” of increase in function f

Chain Rule:

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

$$(3x+2)^2$$

$$g(x): 3x+2$$

$$\begin{aligned} \frac{\partial f(g(x))}{\partial x} &= 2g(x) \cdot \frac{\partial g(x)}{\partial x} \\ &= 2(3x+2) \cdot 3 \end{aligned}$$

Deep network notation

Single data point (row vector) $\left[\leftarrow * \rightarrow \right]$

Input: $X = Z_0 \in \mathbb{R}^{1 \times d_0}$

Intermediate layer: $Z_{i+1} = \sigma(Z_i W_i) \in \mathbb{R}^{1 \times d_{i+1}}$

$Z_i \in \mathbb{R}^{1 \times d_i}$ $W_i \in \mathbb{R}^{d_i \times d_{i+1}}$

Output: $h_\theta(x) = Z_{L+1} \in \mathbb{R}^{1 \times d_{L+1}}$

Loss: $l(Z_{L+1}, y)$: cross entropy (Z_{L+1}, y)

$$z_{i+1} = \sigma(z_i w_i)$$

$$h_\theta(x) = z_L$$

$$\frac{\partial l(h_\theta(x), y)}{\partial w_i} = \frac{\partial l(z_{L+1}, y)}{\partial w_i}$$

chain
rule

$$= \frac{\partial l}{\partial z_{L+1}} \cdot \frac{\partial z_{L+1}}{\partial z_L} \cdot \dots \cdot \frac{\partial z_{i+2}}{\partial z_{i+1}} \cdot \frac{\partial z_{i+1}}{\partial w_i}$$

$$\frac{\partial l(z_{L+1}, y)}{\partial w_{i-1}} = \frac{\partial l}{\partial z_{L+1}} \cdot \frac{\partial z_{L+1}}{\partial z_L} \cdots \frac{\partial z_{i+2}}{\partial z_{i+1}} \cdot \frac{\partial z_{i+1}}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{i-1}}$$

$$\frac{\partial l(z_{L+1}, y)}{\partial w_i} = \frac{\partial l}{\partial z_{L+1}} \cdot \frac{\partial z_{L+1}}{\partial z_L} \cdots \frac{\partial z_{i+2}}{\partial z_{i+1}} \cdot \frac{\partial z_{i+1}}{\partial w_i}$$

avoid repeated computations

$$\frac{\partial l(z_{L+1}, y)}{\partial z_{i+1}} = \delta_{i+1}$$

$$G_i = \frac{\partial l(Z_{L+1}, y)}{\partial W_i} \in \mathbb{R}^{d_i \times d_{i+1}}$$

$$\begin{aligned} Z_i &\in \mathbb{R}^{1 \times d_i} \\ Z_{i+1} &\in \mathbb{R}^{1 \times d_{i+1}} \\ W_i &\in \mathbb{R}^{d_i \times d_{i+1}} \end{aligned}$$

$$\frac{\partial l(Z_{L+1}, y)}{\partial Z_i} = \frac{\partial l(Z_{L+1}, y)}{\partial Z_{i+1}} \cdot \frac{\partial Z_{i+1}}{\partial Z_i} \quad Z_{i+1} = \sigma(Z_i W_i)$$

$\swarrow \quad \searrow$
 $1 \times d_{i+1} \quad d_{i+1} \times d_i$
 $\in \mathbb{R} \quad \in \mathbb{R}$

$$G_i = \left(G_{i+1} \odot \underbrace{\sigma'(Z_i W_i)}_{\in \mathbb{R}^{1 \times d_{i+1}}} \right) W_i^T$$

$\mathbb{R}^{1 \times d_{i+1}} \quad \mathbb{R}^{d_{i+1} \times d_i}$

$$\in \mathbb{R}^{d_{i+1} \times d_i}$$

Backward iteration

⊙: element wise product

$$G_i = \frac{\partial l(Z_{L+1}, y)}{\partial Z_i}$$

$$G_i = G_{i+1} \odot \sigma'(Z_i W_i) \cdot W_i^T$$

$\frac{\partial l(Z_{L+1}, y)}{\partial W_i}$: a bit tricky to derive from first principles

"Hack": write down different terms from chain rule and make dimensions match

$$\begin{aligned} \frac{\partial l(Z_{L+1}, y)}{\partial W_i} &= \frac{\partial l(Z_{L+1}, y)}{\partial Z_{i+1}} \cdot \frac{\partial Z_{i+1}}{\partial W_i} \\ &= \underset{\substack{\mathbb{R} \xrightarrow{d_i \times d_{i+1}}}}{\downarrow} G_{i+1} \underset{\substack{\mathbb{R} \xrightarrow{1 \times d_{i+1}}}}{\downarrow} \sigma'(Z_i W_i) \underset{\substack{\mathbb{R} \xrightarrow{1 \times d_i}}}{\downarrow} Z_i^T \rightarrow \mathbb{R} \end{aligned}$$

$$\frac{\partial \ell}{\partial w_i} = \underbrace{(z_i)^T}_{\substack{\downarrow d_i \times 1 \\ \mathbb{R}}} \underbrace{(G_{i+1} \odot \sigma'(z_i w_i))}_{\substack{\downarrow \mathbb{R}^{1 \times d_{i+1}} \\ \mathbb{R}}}$$

Batched : m samples at a time independently
 version (talk more in Lec 9)

$$z_i \in \mathbb{R}^{m \times d_i}$$

$$G_i \in \mathbb{R}^{m \times d_i}$$

$$\begin{bmatrix} \leftarrow z_i^{(1)} \rightarrow \\ \leftarrow z_i^{(2)} \rightarrow \\ \vdots \\ \leftarrow z_i^{(m)} \rightarrow \end{bmatrix}$$

$$G_i = G_{i+1} \odot \sigma'(z_i w_i) w_i^T$$

$\downarrow m \times d_{i+1} \quad \downarrow m \times d_{i+1} \quad \downarrow d_{i+1} \times d_i$
 $\mathbb{R} \quad \mathbb{R} \quad \mathbb{R}$

$$\frac{\partial \ell}{\partial w_i} = \underbrace{(z_i)^T}_{\substack{\downarrow d_i \times m \\ \mathbb{R}}} \underbrace{(G_{i+1} \odot \sigma'(z_i w_i))}_{\substack{\downarrow m \times d_{i+1} \\ \mathbb{R}}}$$

$\downarrow d_i \times d_{i+1}$
 \mathbb{R}

Terminology

Z_i = intermediate values in the network

Z_1 = input

Z_L = output

aka activations

$$Z_i = \sigma(W_{i-1}^T Z_{i-1})$$

activation function

W_i : "weights" & "parameters"

have to be estimated from data

learning process: find best W_i that fit data

$$x^{(i)} \rightarrow Z_L^{(i)}$$

$$\text{loss}(Z_L^{(i)}, y^{(i)})$$

Backpropagation algorithm

- Initialize $Z_1 = X$

- For $i = 1 \dots L$

$$Z_{i+1} = \sigma(Z_i W_i)$$

Forward pass

G_L : gradient of loss function (cross entropy)

- For $i = L \dots 1$

$$G_i = G_{i+1} \odot \sigma'(Z_i W_i) W_{i+1}^T$$

Backward pass

$$\frac{\partial l}{\partial W_i} = Z_i^T G_{i+1} \odot \sigma'(W_i^T Z_i)$$

each element

Training algo:

$$W^{(t)} = W^{(t-1)} - \eta \left[\nabla_w l \right]$$

$[W_1, W_2, \dots, W_L]$

Tr Loss: $\sum_{i=1}^n \text{loss}(x^{(i)}, y^{(i)})$

$\nabla \text{Tr Loss} = \nabla \sum = \sum \nabla$

Stochastic: gradient of a batch of samples

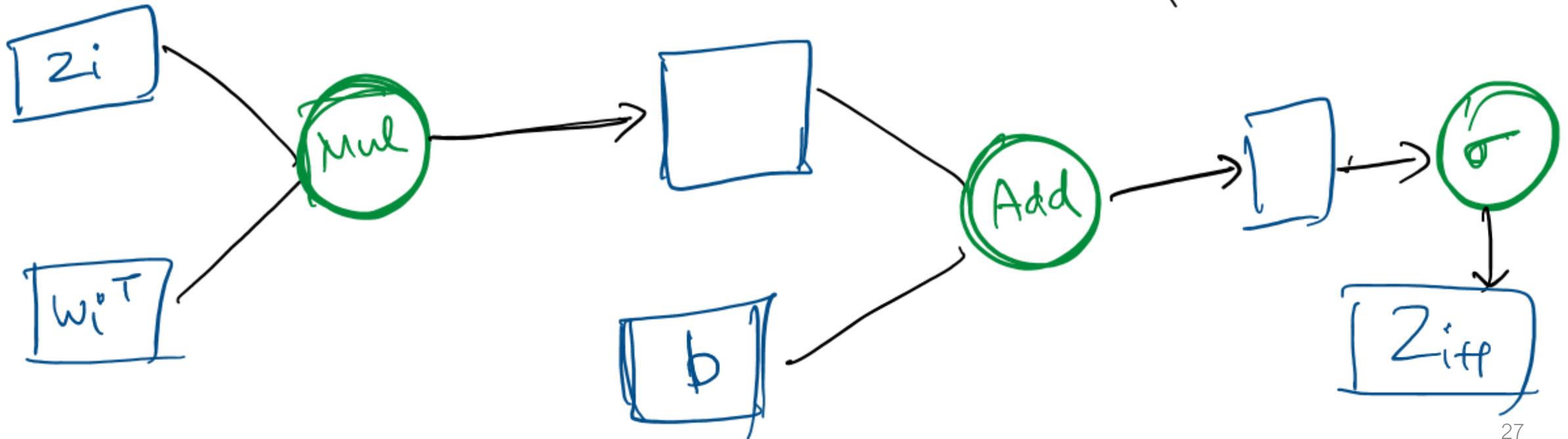
Computation graph

\square : variables

\bigcirc : functions

Directed Acyclic Graph

~~$z_{iH} = \sigma(w_i^T z_i)$~~
 $z_{iH} = \sigma(w_i^T z_i + b_i)$



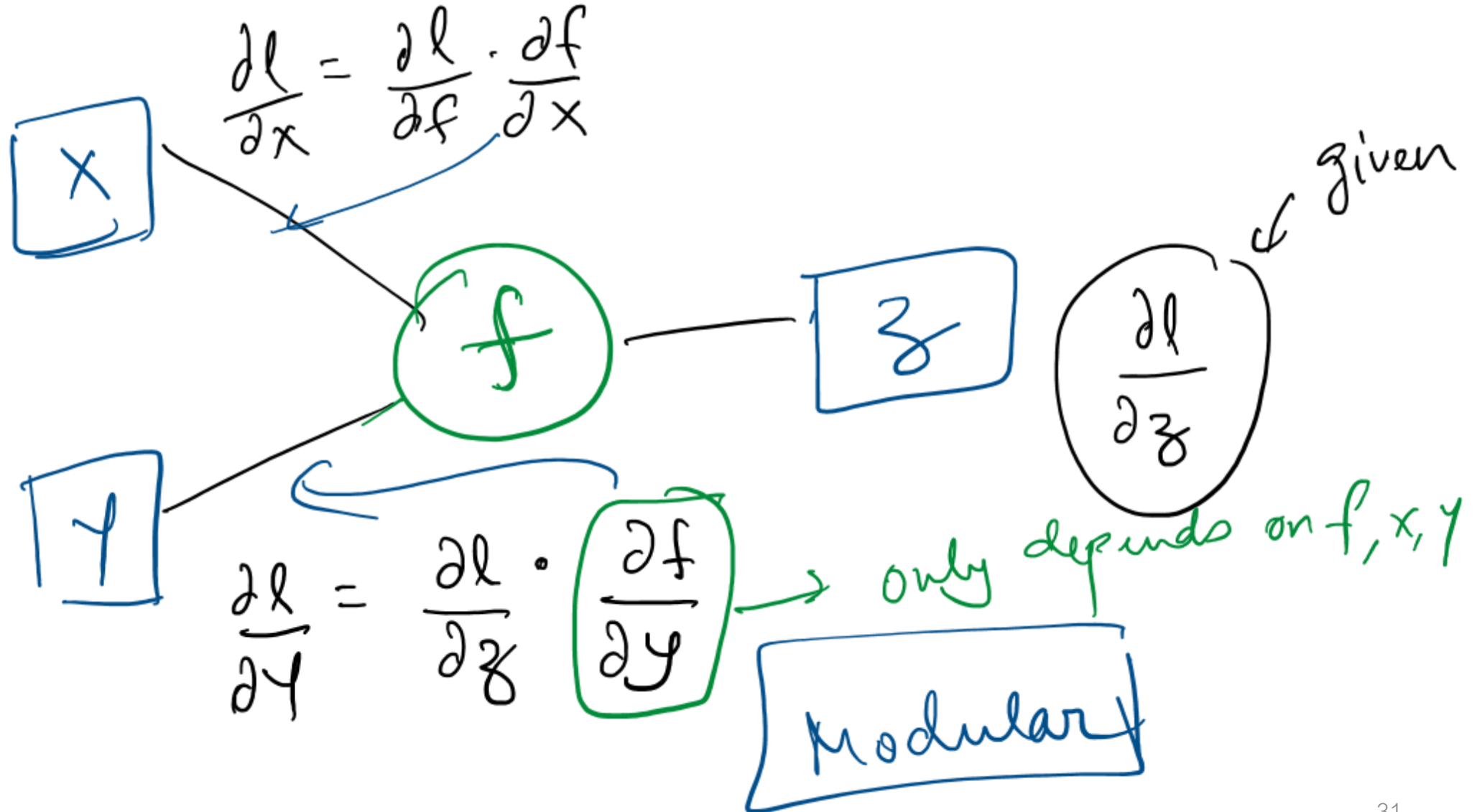
Computation graph

Directed Acyclic Graph to represent the functions computed

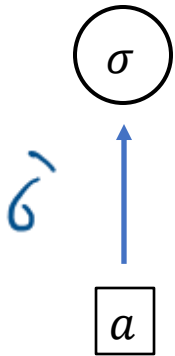
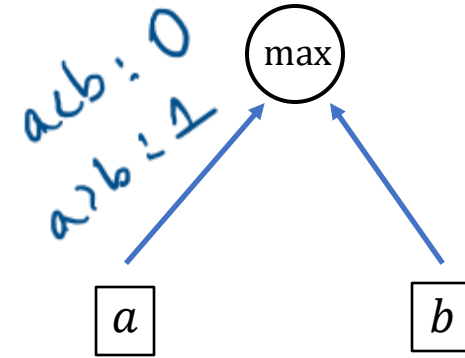
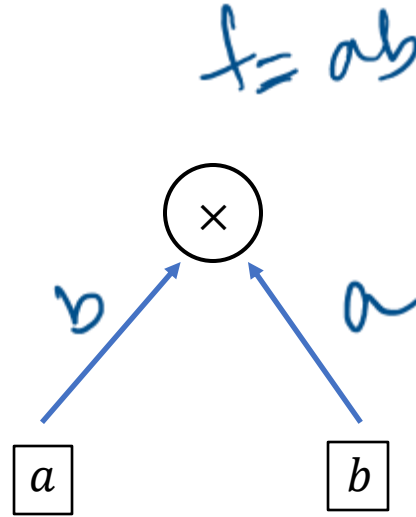
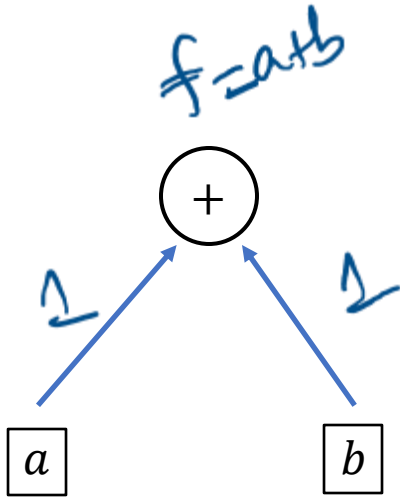
- Root node has the final expression, intermediate nodes have subexpressions
- Convenient to compute gradients, used in popular frameworks like pytorch and tensorflow

Computation graph example

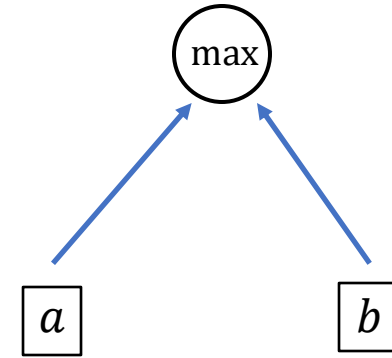
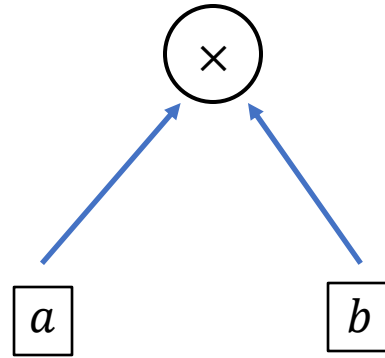
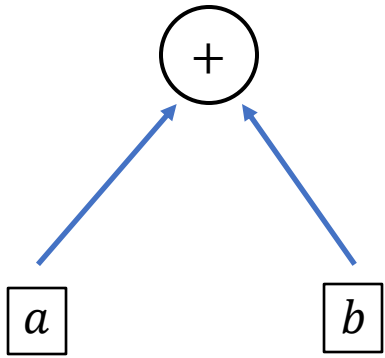
Computation graph chain rule



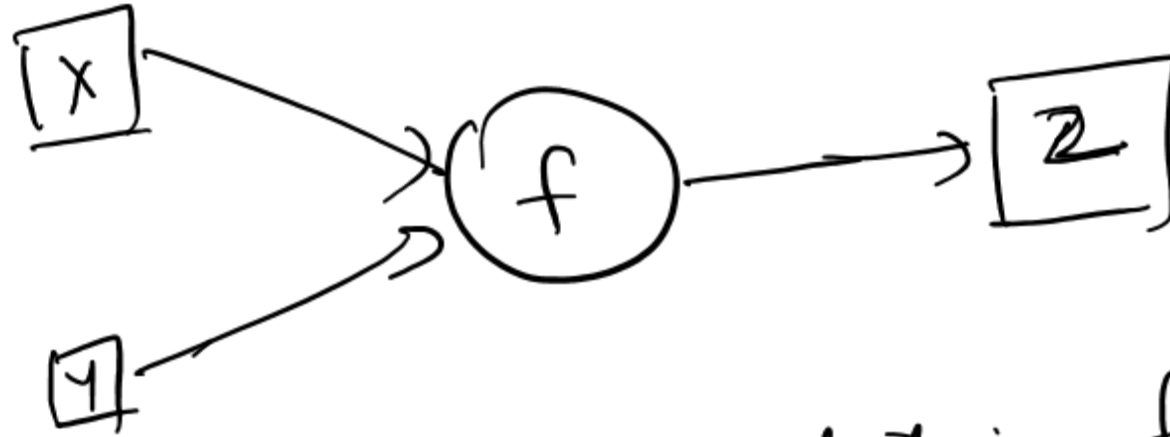
Computation graph: gradients along edges



Computation graph



Auto diff



Forward pass: implementation $f(x, y)$

Backward pass:

Given $\frac{\partial l}{\partial z}$

vector

$$\frac{\partial l}{\partial z} \cdot \frac{\partial z}{\partial x},$$

$$\frac{\partial l}{\partial z} \cdot \frac{\partial z}{\partial y} \dots$$

in general: inputs: $[x_1, x_2 \dots x_n]$

output: $[z_1 \dots z_d]$

$$\begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_2}{\partial x_1} & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ \dots & \dots & \dots & \frac{\partial z_d}{\partial x_n} \end{bmatrix}$$

→ high dim

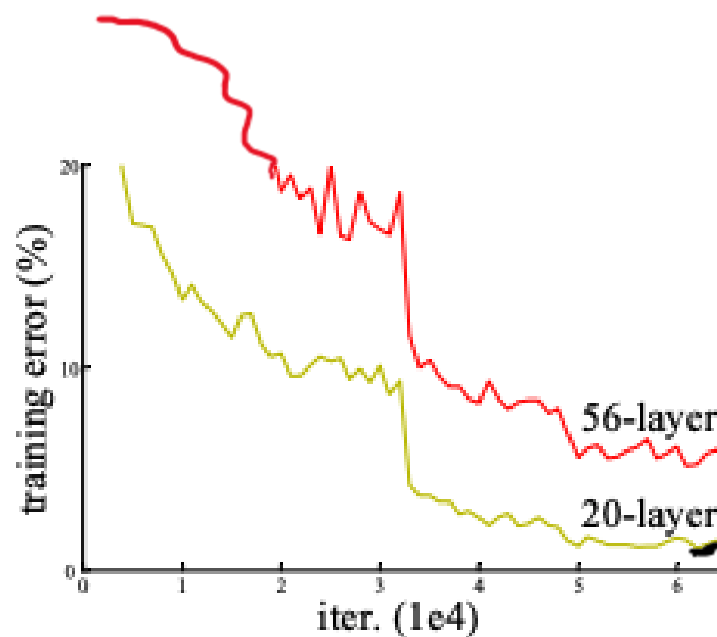
$$f: \mathbb{R}^n \rightarrow \mathbb{R}^d$$

Jacobian

Multi-layer perceptron recap

- feedforward networks
- "more layers" are better
- Perceptron by Minsky \equiv AI winter

Piazza poll question



optimization
issue

7 56-layer network

} identity

56 layer

20 layer

20
layer

Piazza poll question

→ \exists a 56 layer network with low train error

[identity]

20
layer ["good"]

Optimization issue

vanishing gradient:

each "edge" \equiv a multiplication
and it could make the gradient
really small

exploding

Residual connection

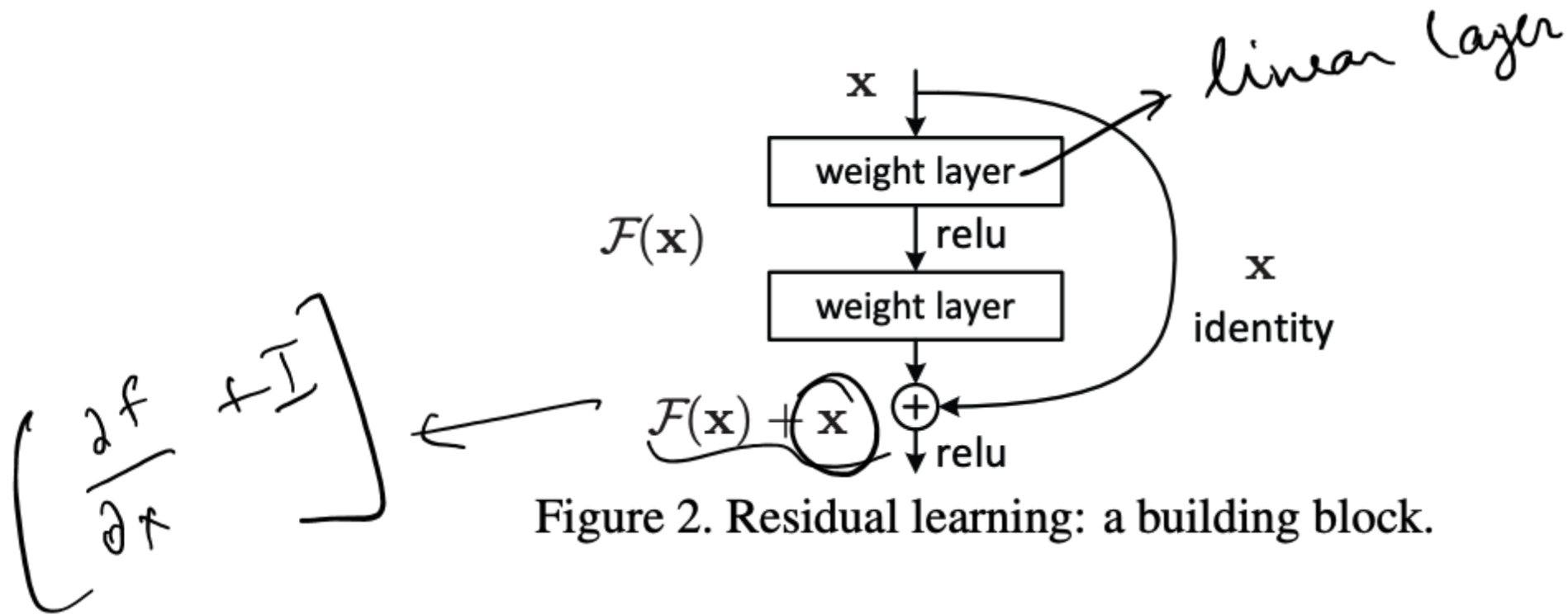


Figure 2. Residual learning: a building block.

Residual connection

Two interpretations

1. Easier to preserve identity / good features
2. Addresses vanishing gradients due to “shortcut” connections

Normalization

→ Z_i s also change a lot : compounding across layers

generic form:

$$\hat{Z}_i = \left(\frac{Z_i - \mathbb{E}[Z_i]}{\sqrt{\text{var}(Z_i) + \epsilon}} \right) \cdot \cancel{\gamma} + \cancel{\lambda}$$

$Z_i \in \mathbb{R}^d$

to avoid divide by zero

\Downarrow
 Z_i are activation
of a layer \equiv layer norm

$\gamma, \lambda \in \mathbb{R}^d$

λ : bias
 γ : scale

} learnable

The canonical architecture: transformers

