

S&DS 365 / 565
Data Mining and Machine Learning

Yale

For today

Recap, go to notebook

Neural networks: Fully-connected feed-forward and backpropagation

Non-linear methods

Feature mappings

Advanced optimization

Recap

Optimization: gradient descent, newton's method, stochastic gradient descent

SGD Goal:

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

SGD Algo:

$$\theta_k = \theta_{k-1} - \eta_{k-1} \nabla f_J(\theta_{k-1})$$

$$J \sim [n]$$

Recap

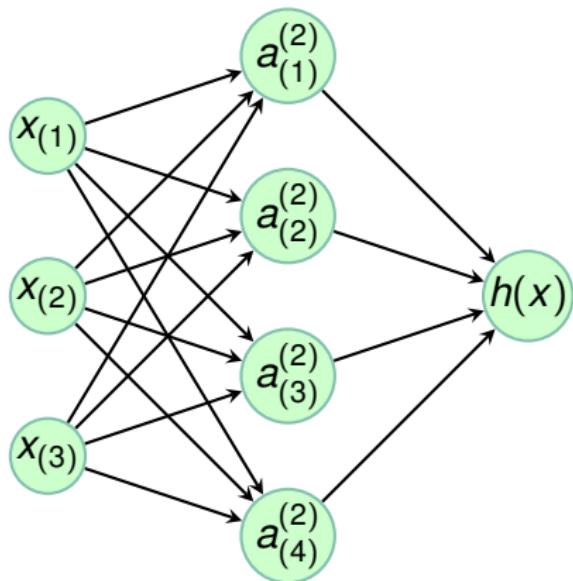
Example: $\{x_i, y_i\}_{i=1}^n$ with logistic

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n [\log(1 + \exp(x_i^T \theta)) - y_i x_i^T \theta]$$

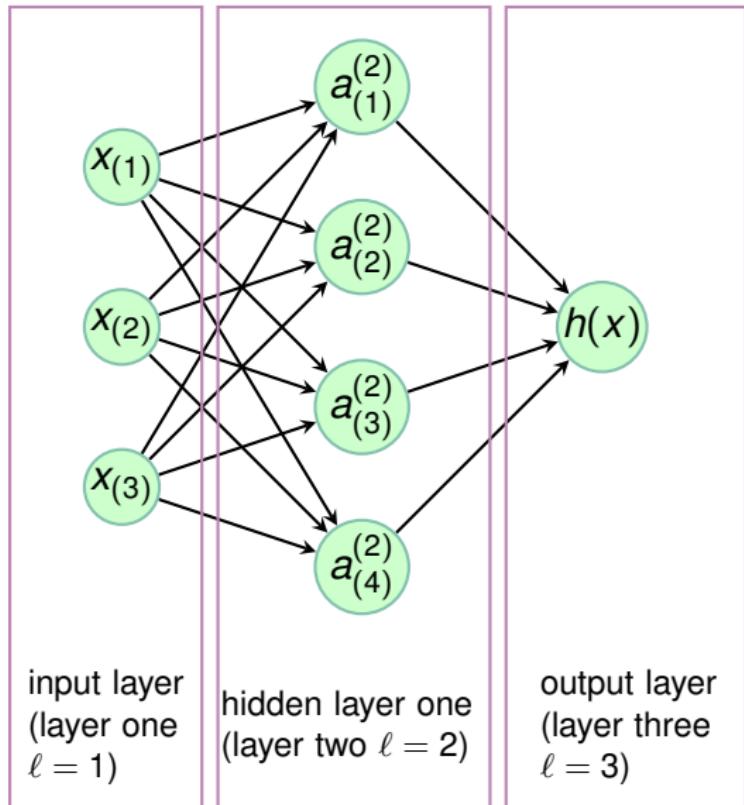
$$f_i(\theta) = \log(1 + \exp(x_i^T \theta)) - y_i x_i^T \theta$$

$$\theta_{k+1} = \theta_k - \eta_k \left[\frac{\exp(x_J^T \theta_k)}{1 + \exp(x_J^T \theta_k)} - y_J \right] x_J$$

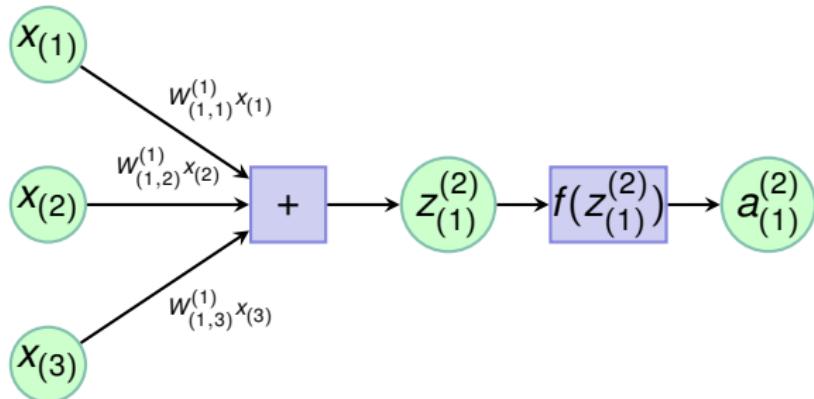
Neural network



Neural network



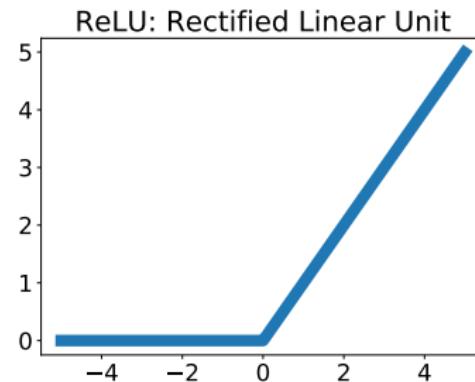
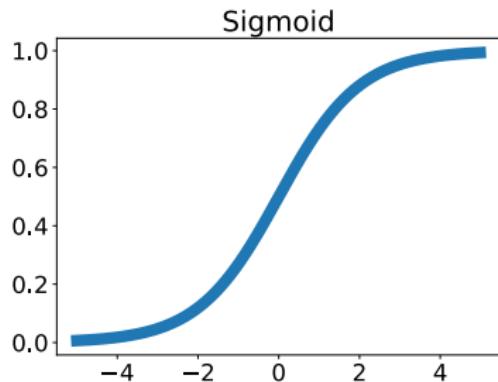
Neural network



Let $w_1^1 = W_{(1,:)}^{(1)} \in \mathbb{R}^3$

$$a_{(1)}^2 = f(x^T w_1^1) = f(\langle x, w_1^1 \rangle)$$

f is activation function



Neural network: three layers $L = 3$

Weights of layers

$$W^{(1)} \in \mathbb{R}^{d_2 \times d}$$

$$W^{(2)} \in \mathbb{R}^{1 \times d_2}$$

$$\begin{array}{lll} x \in \mathbb{R}^d & f: \mathbb{R}^d \rightarrow \mathbb{R}^{d_2} \\ b^{(1)} \in \mathbb{R}^{d_2} & b^{(2)} \in \mathbb{R} & h: \mathbb{R}^d \rightarrow \mathbb{R} \end{array}$$

$$h(x) = W^{(2)}f(W^{(1)}x + b^{(1)}) + b^{(2)}$$

Note

$$[f(v)]_{(i)} = f(v_{(i)})$$

If f non-linear $\in \mathbb{R}^{1 \times d}$: e.g. $f(v) = v$

$$\Rightarrow h(x) = \underbrace{W^{(2)}W^{(1)}}_{\text{Linear function!}} x + (W^{(2)}b^{(1)} + b^{(2)})$$

Neural Network simple example $L = 2$

$$W^{(1)} \in \mathbb{R}^{1 \times d}$$

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

$$h: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$h(x) = W^{(1)}x$$

Neural network: four layers $L = 4$

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

$$W^{(1)} \in \mathbb{R}^{d_2 \times d}$$

$$W^{(2)} \in \mathbb{R}^{d_3 \times d_2}$$

$$W^{(3)} \in \mathbb{R}^{1 \times d_3}$$

$$b^{(1)} \in \mathbb{R}^{d_2}$$

$$b^{(2)} \in \mathbb{R}^{d_3}$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^{d_2}$$

$$f: \quad \quad \quad \rightarrow \mathbb{R}^{d_3}$$

$$h: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$h(x) = f(W^{(3)}a^{(3)})$$

$$a^{(3)} = f(W^{(2)}a^{(2)} + b^{(2)})$$

$$a^{(2)} = f(W^{(1)}x + b^{(1)})$$

$$x$$

④

③

②

①

Neural network: more general outputs

$$W^{(3)} \in \mathbb{R}^{d_4 \times d_3}$$

vector

$$h(x) = f(W^{(3)} a^{(3)})$$

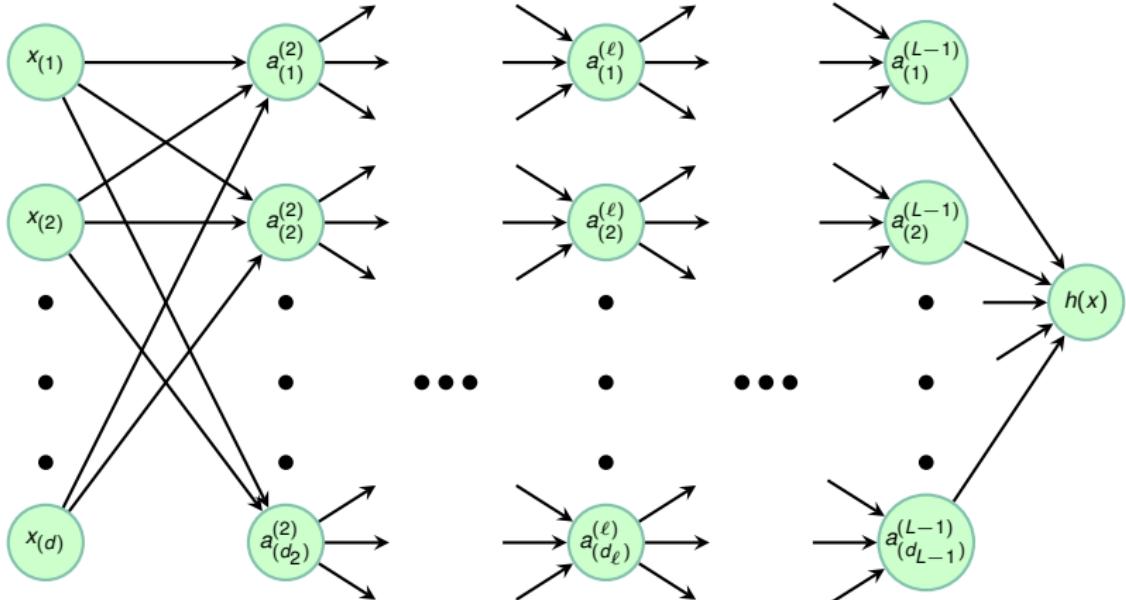
$$a^{(3)} = f(W^{(2)} a^{(2)} + b^{(2)})$$

$$a^{(2)} = f(W^{(1)} x + b^{(1)})$$

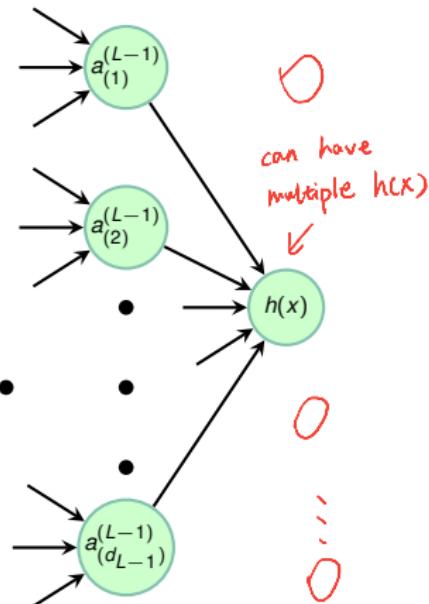
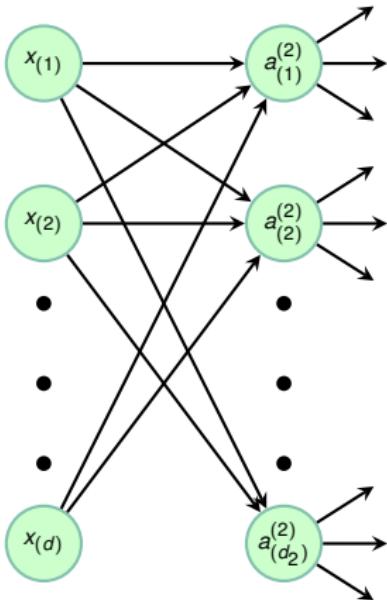
Neural Networks

Fully-connected feed-forward

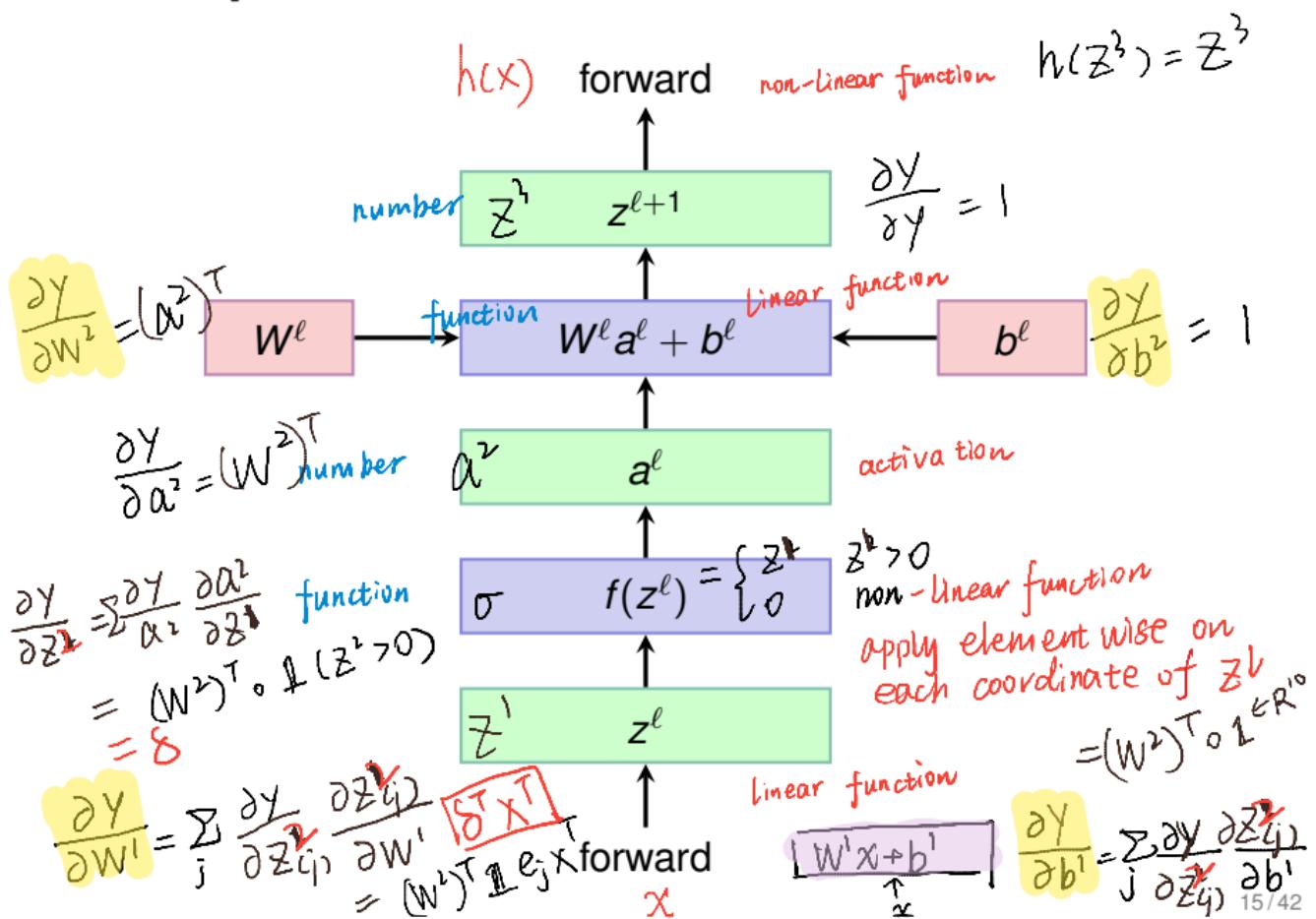
Neural network: L layers



Neural network



Forward pass



Neural network: math

Parameters: $W^{(\ell)} \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$, $b^{(\ell)} \in \mathbb{R}^{d_{\ell+1}}$

Input: $x \in \mathbb{R}^d$ and $d = d_1$

Output: $h(x) = z^{(L)}$ and $d_L \geq 1$

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)}) \quad \text{f: nonlinear function}$$

$$a^{(k)} = f(z^{(k)})$$

$$z^{(k+1)} = W^{(k)}a^{(k)} + b^{(k)}$$

we also train neural network by gradient descent, so need to compute gradient

How do we compute gradients?

$w^{(l)}, b^{(l)}$ are params
though w and b may not in function h
 h still depend on w and b

subscripts W and b included above

generally will drop them

$$\left\{ \begin{array}{l} \frac{\partial h_{W,b}(x)}{W^{(l)}} \\ \frac{\partial h_{W,b}(x)}{b^{(l)}} \end{array} \right.$$

loss a choice: mean squared loss

$$\Rightarrow L(W, b) = \frac{1}{2n} \sum_{i=1}^n (h(x_i) - y_i)^2$$

$$\Rightarrow \frac{\partial L}{\partial W} = \frac{1}{n} \sum_{i=1}^n (h(x_i) - y_i) \frac{\partial h(x)}{\partial W}$$

对矩阵求梯度
矩阵

Recall chain rule

$g : \mathbb{R}^p \mapsto \mathbb{R}$ and $f : \mathbb{R}^d \mapsto \mathbb{R}^p$

intermediary variable ✓

Let $v = f(\beta)$ and $y = g(v)$

multivariable vector
change β of i change every coordinate of v of j

one form

$$\frac{\partial y}{\partial \beta(i)} = \sum_{j=1}^p \frac{\partial y}{\partial v(j)} \frac{\partial v(j)}{\partial \beta(i)}$$

gradient

$$= \sum_{j=1}^p \underbrace{\left[\nabla g(v) \Big|_{v=f(\beta)} \right]_{(j)}}_{\text{forward pass}} \frac{\partial f(\beta)_{(j)}}{\partial \beta(i)}$$

Jacobian of f

another form

$$\frac{\partial g(f(\beta))}{\partial \beta(i)} = \sum_{j=1}^p \frac{\partial g(f(\beta))}{\partial [f(\beta)]_{(j)}} \frac{\partial [f(\beta)]_{(j)}}{\partial \beta(i)}$$

Let $v, w \in \mathbb{R}^d$

Recall: $f(v) = v^T w \Rightarrow \nabla f(v) = \frac{\partial f(v)}{\partial v} = w$

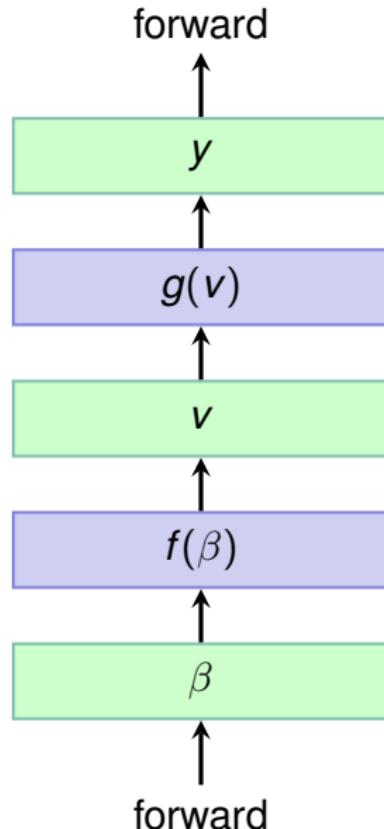
Notation: $z = w \circ v$ point-wise multiplication Hadamard product

$$z_{(i)} = w_{(i)} v_{(i)}$$

standard base vector only jth coordinate is non-zero jth coordinate of v
Let $e_j \in \mathbb{R}^d$ such that $[e_j]_{(i)} = 1(i=j) \Rightarrow e_j^T v = v_{(j)}$

$$v = \sum_j e_j v_{(j)}$$
$$v = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \end{pmatrix} v_{(1)} + \begin{pmatrix} 0 \\ \text{sum} \\ 0 \\ 1 \\ 0 \end{pmatrix} v_{(2)} + \dots$$
$$v \circ w = \sum_j e_j v_{(j)} w_{(j)}$$

Graphically

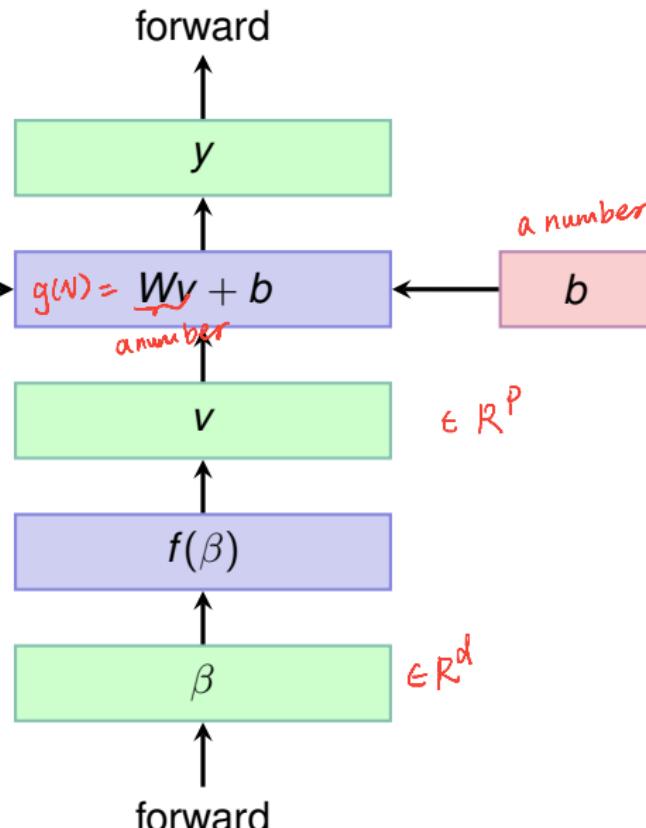


Example

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

$$g: \mathbb{R}^P \rightarrow \mathbb{R}$$

$W \in \mathbb{R}^{1 \times P}$ *a row vector*



Mathematically

$$y = Wf(\beta) + b$$

$$W \in \mathbb{R}^{1 \times p}$$

$$b \in \mathbb{R}$$

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

Gradient

compute

$$y = Wf(\beta) + b$$

Want: $\frac{\partial y}{\partial W}, \frac{\partial y}{\partial b}, \frac{\partial y}{\partial \beta}$

$$g(v) = Wv + b \quad W \text{ is a row vector}$$

$$\nabla_v g(v) = W^T \quad \text{now a column vector}$$

Gradient

$$y = Wf(\beta) + b$$

Want: $\frac{\partial y}{\partial W}$, $\frac{\partial y}{\partial b}$, $\frac{\partial y}{\partial \beta}$

$$g(v) = Wv + b$$

$$\nabla_v g(v) = W^T \text{ now a column vector}$$

meaning of $\frac{\partial y}{\partial \beta(i)}$: if changing i th coordinate of β , how does that change j th coordinate of $f(\beta)$ i.e. how does it change j th coordinate of v , then how does changing j th coordinate of v change i th coordinate of $f(\beta)$ just change by $w_{c(j)}$

$$\begin{aligned}\frac{\partial y}{\partial \beta(i)} &= \sum_{j=1}^p \left[\nabla g(v) \Big|_{v=f(\beta)} \right]_{(j)} \frac{\partial f(\beta)_{(j)}}{\partial \beta(i)} \\ &= \sum_{j=1}^p W_{(j)} \frac{\partial f(\beta)_{(j)}}{\partial \beta(i)}\end{aligned}$$

Gradient

$$y = Wf(\beta) + b$$

Want: $\frac{\partial y}{\partial W}$, $\frac{\partial y}{\partial b}$, $\frac{\partial y}{\partial \beta}$

$$g(v) = Wv + b$$

$$\nabla_v g(v) = W^T \quad \text{now a column vector}$$

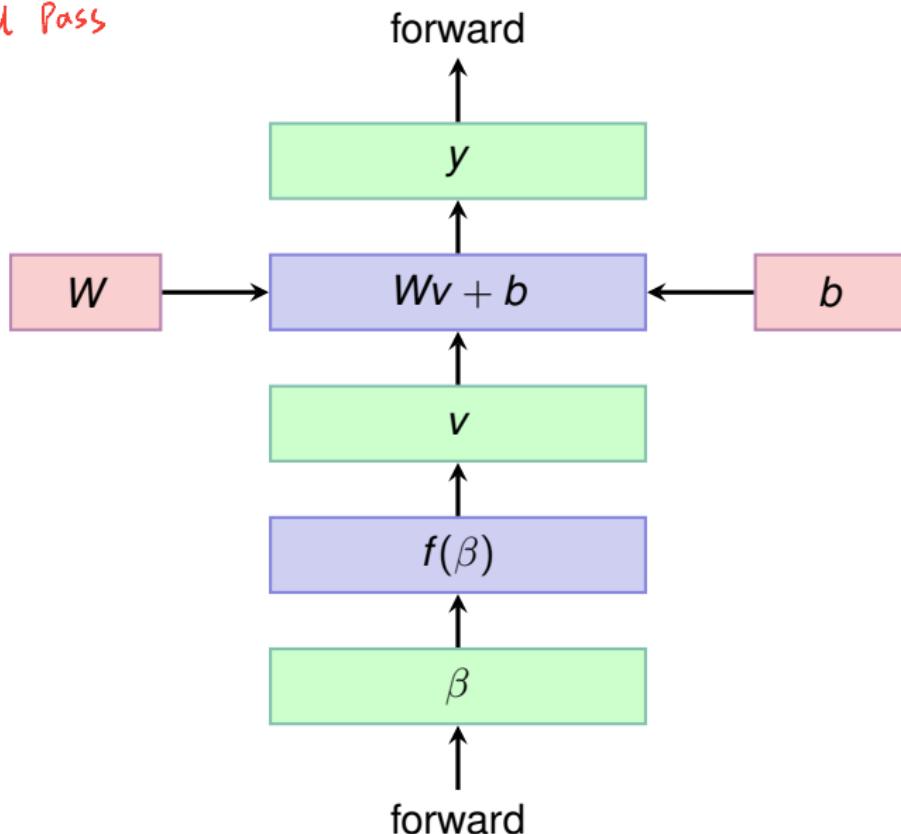
easy to find

$$\left\{ \begin{array}{l} \frac{\partial y}{\partial W} = f(\beta)^T \\ \frac{\partial y}{\partial b} = 1 \end{array} \right.$$

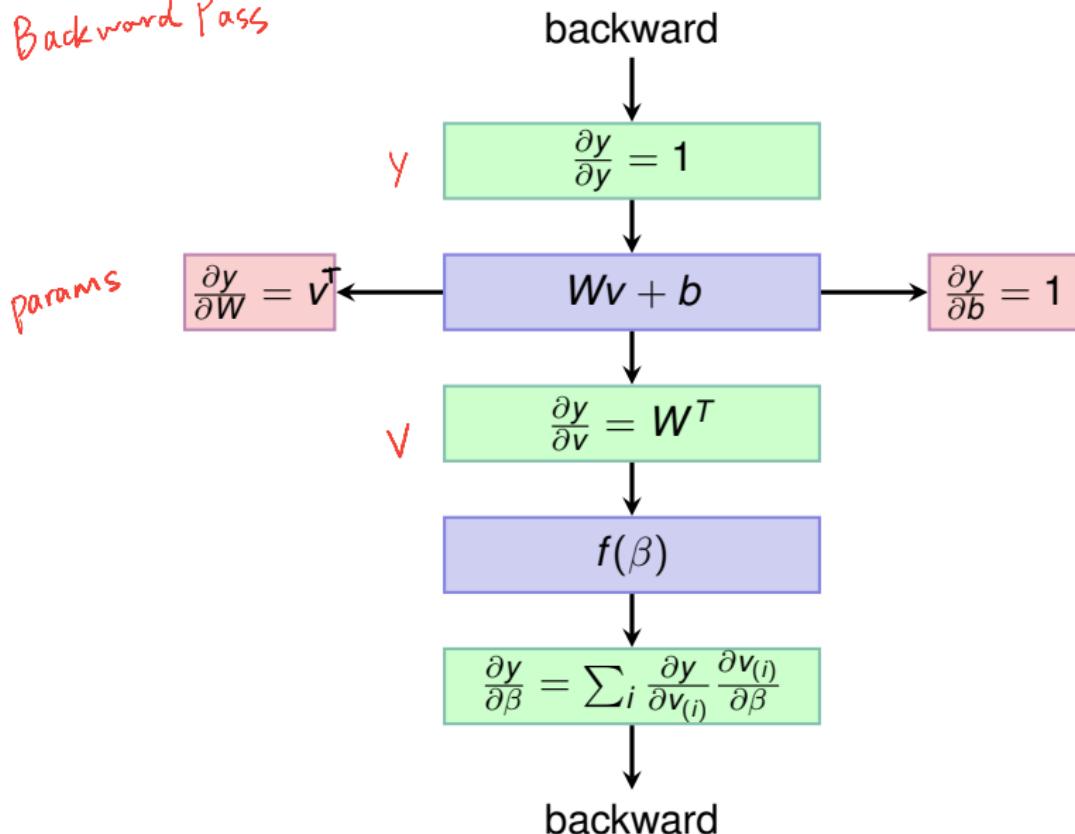
gradient should be the same shape as W

a row vector

Forward Pass



Backward Pass

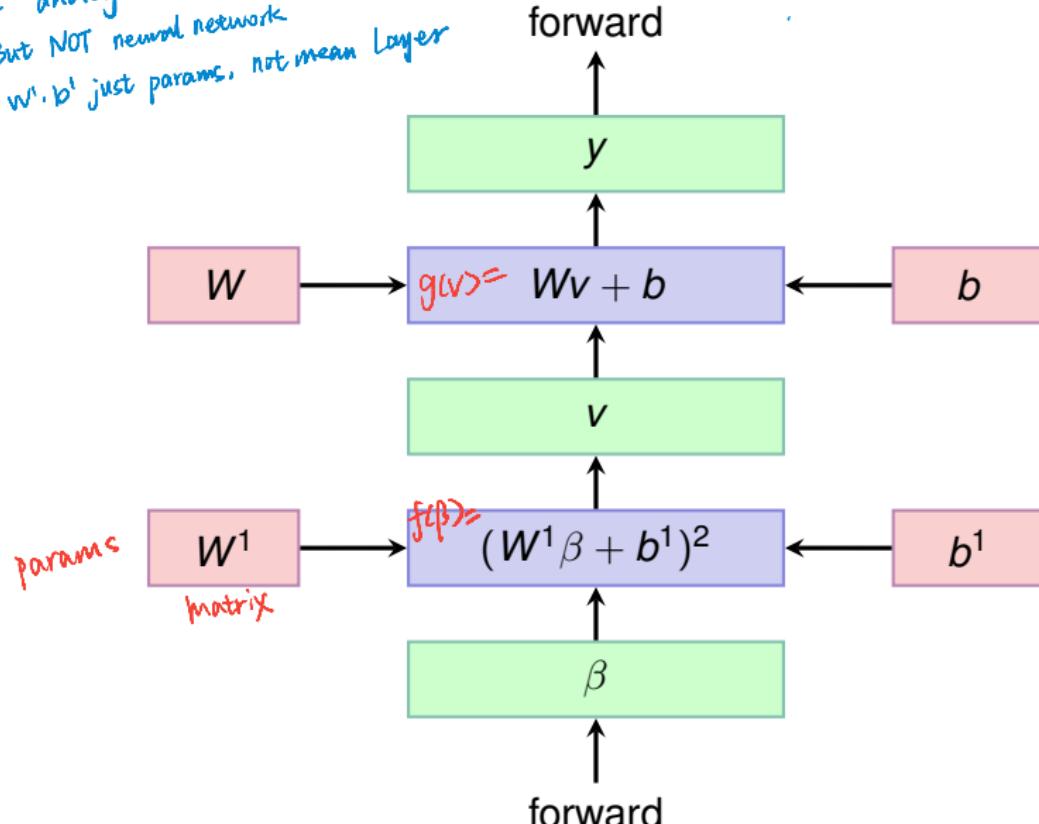


More concrete example

a analogue to Neural Network

But NOT neural network

w, b¹ just params, not mean layer



Mathematically

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^p \quad g: \mathbb{R}^p \rightarrow \mathbb{R}$$

$$W \in \mathbb{R}^{1 \times p} \quad \text{a row vector}$$

$$b \in \mathbb{R}$$

$$W^1 \in \mathbb{R}^{p \times d} \quad \text{a matrix}$$

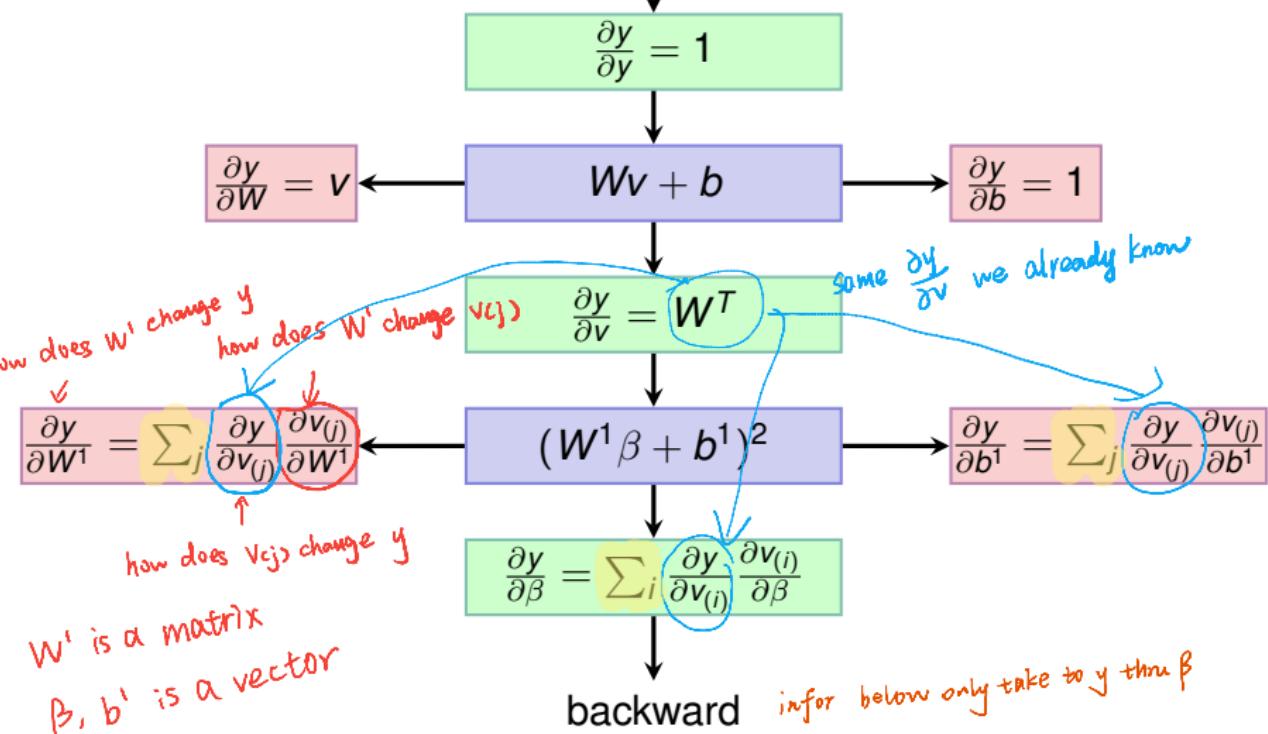
$$\beta \in \mathbb{R}^d \quad b^1 \in \mathbb{R}^p$$

$$y = W((W^1\beta + b^1)^2) + b$$

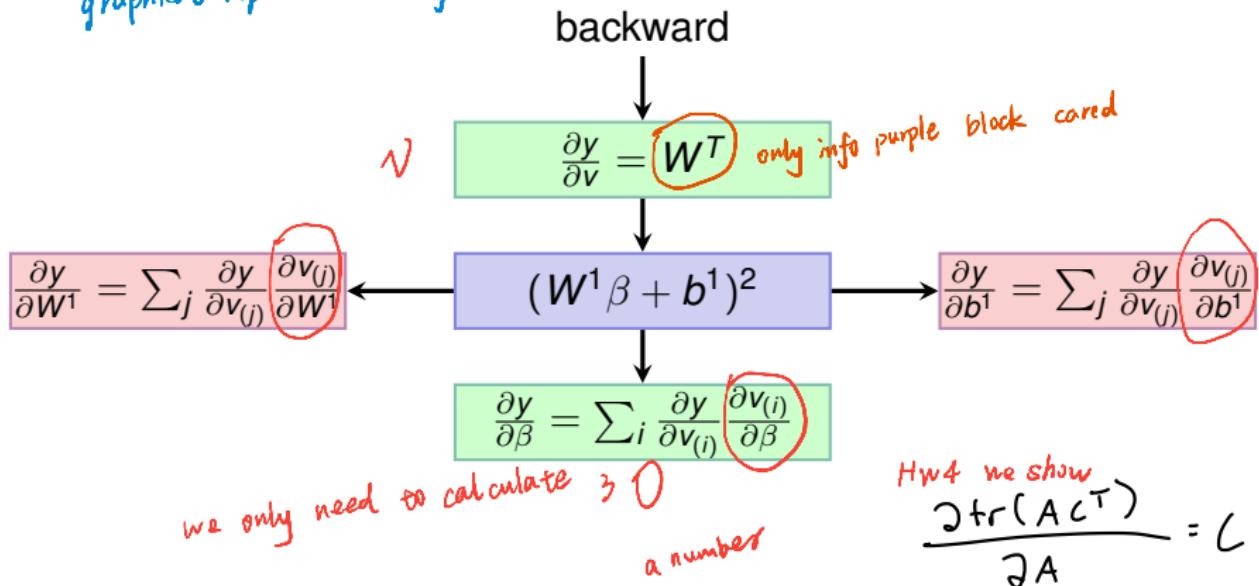
$$[(W^1\beta + b^1)^2]_{(i)} = [W^1\beta + b^1]_{(i)}^2$$

Chain rule of backward pass

backward



graphical representation of chain rule



since $e_j^T W^1 \beta$ is a real number
 $\text{trace}(a) = a \quad a \in \mathbb{R}$

$$e_j^T W^1 \beta = \text{trace}(e_j^T W^1 \beta)$$

$$v_{(j)} = (\underbrace{e_j^T W^1 \beta + b_{(j)}^1}_\text{ej ER^d \beta ER^d})^2$$

$$= (\text{trace}(W^1 \beta e_j^T) + b_{(j)}^1)^2$$

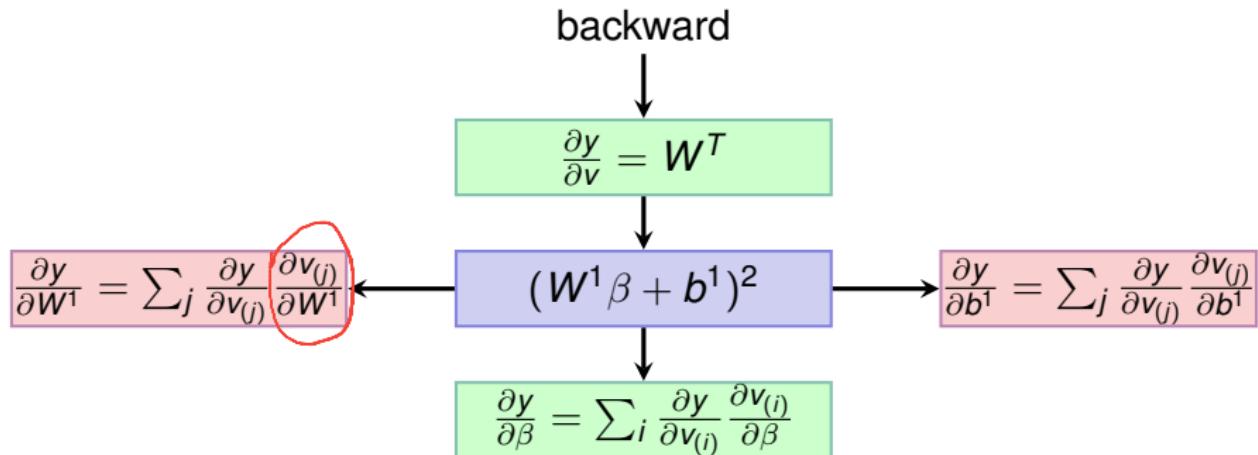
$$\begin{matrix} \text{jth row} \\ \text{of E} \end{matrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} [P_1 P_2 \dots P_j \dots P_d] = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

$$\frac{\partial \text{tr}(A C^T)}{\partial A} = C$$

$$\Rightarrow \frac{\partial \text{tr}(W^1 \beta e_j^T)}{\partial W^1}$$

$$= e_j \beta^T$$

a matrix $E \in \mathbb{R}^{d \times d}$
 only jth row is no zero = β^T



Now compute $\frac{\partial v_{(j)}}{\partial W^1}$

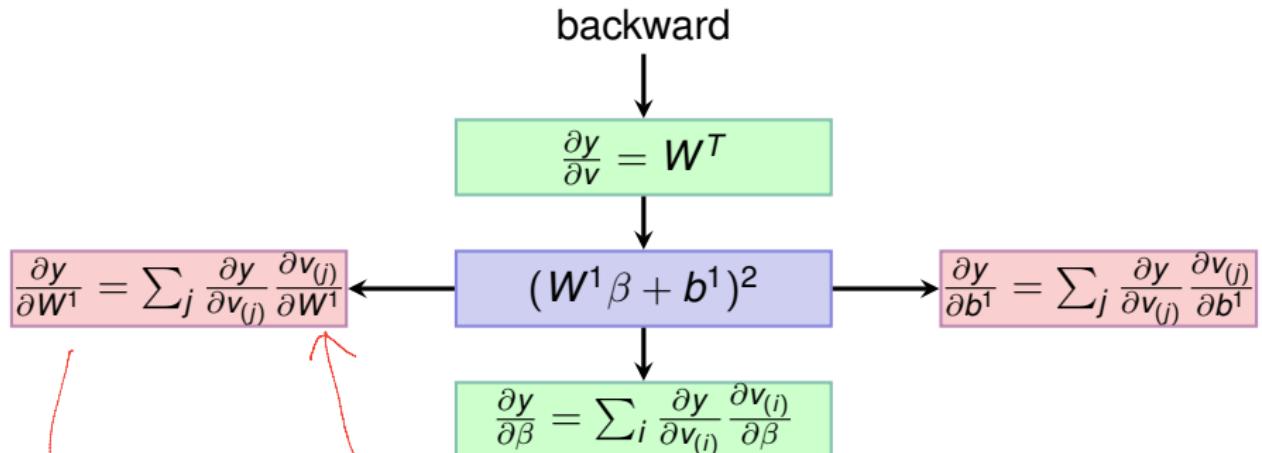
$$\begin{aligned}
 v_{(j)} &= (e_j^T W^1 \beta + b_{(j)}^1)^2 \\
 &= (\text{trace}(W^1 \beta e_j^T) + b_{(j)}^1)^2 \quad = (z_{(j)})^2
 \end{aligned}$$

define

$$z = W^1 \beta + b^1 \Rightarrow z_{(j)} = \text{trace}(W^1 \beta e_j^T) + b_{(j)}^1$$

$$\begin{aligned}
 \frac{\partial v_{(j)}}{\partial W^1} &= 2z_{(j)}(e_j \beta^T) \\
 &= \frac{\partial(z_{(j)})^2}{\partial z_{(j)}} \frac{\partial z_{(j)}}{\partial W^1} \quad \text{chain rule}
 \end{aligned}$$

backward



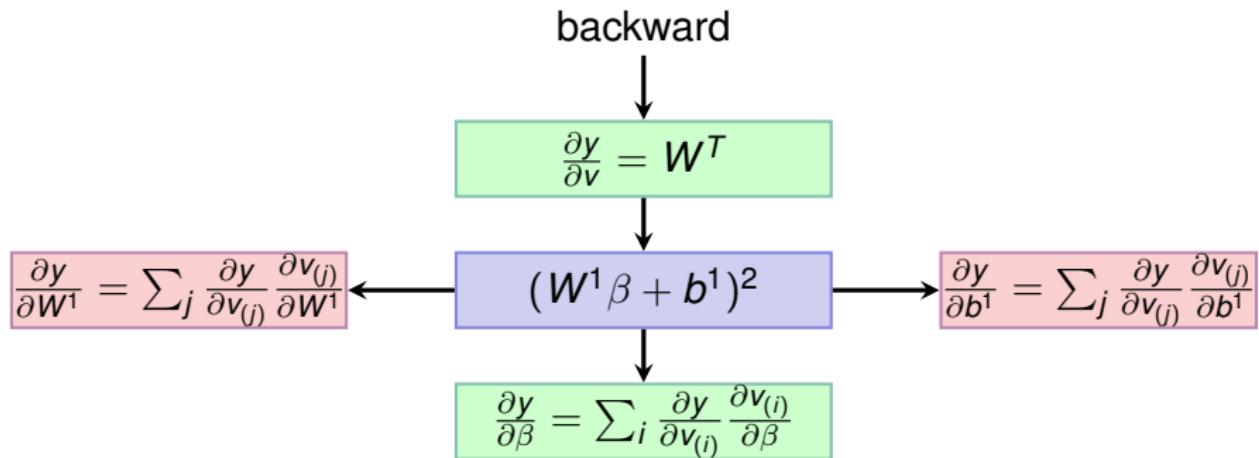
$$z = W^1 \beta + b^1 \Rightarrow z_{(j)} = \text{trace}(W^1 \beta e_j^T) + b_{(j)}$$

plug in

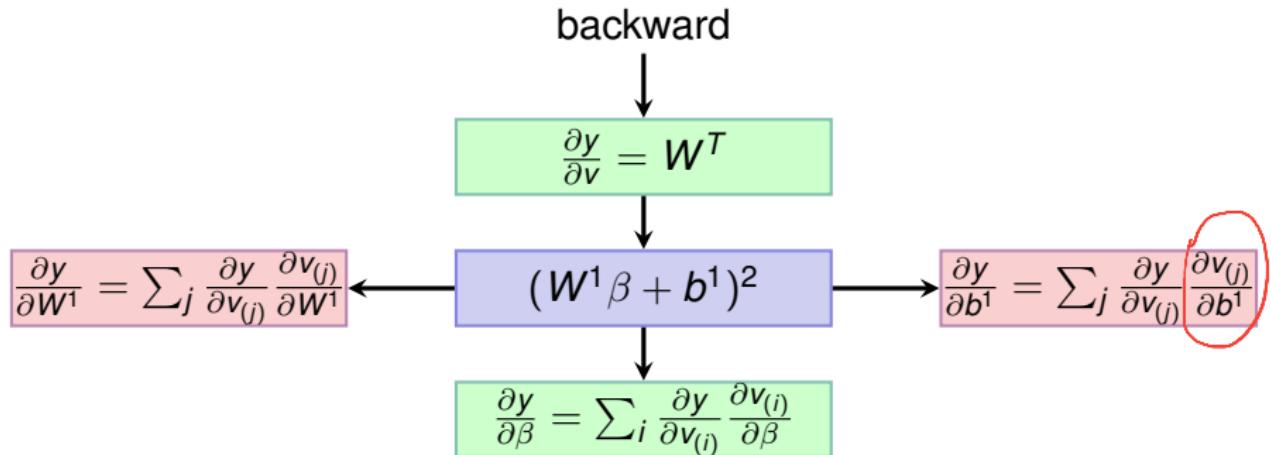
$$\frac{\partial v_{(j)}}{\partial W^1} = 2z_{(j)} e_j \beta^T$$

$$\frac{\partial y}{\partial W^1} = \sum_j 2 \underbrace{W_{(j)} z_{(j)}}_{\frac{\partial y}{\partial v_{(j)}}} e_j \beta^T = 2 \underbrace{(W^T \circ z)}_{\text{recall slide 20}} \beta^T$$

independent of sum



$$v_{(j)} = (\mathbf{e}_j^T W^1 \beta + b_{(j)}^1)^2$$



Now compute $\frac{\partial v_{(j)}}{\partial b^1}$

define

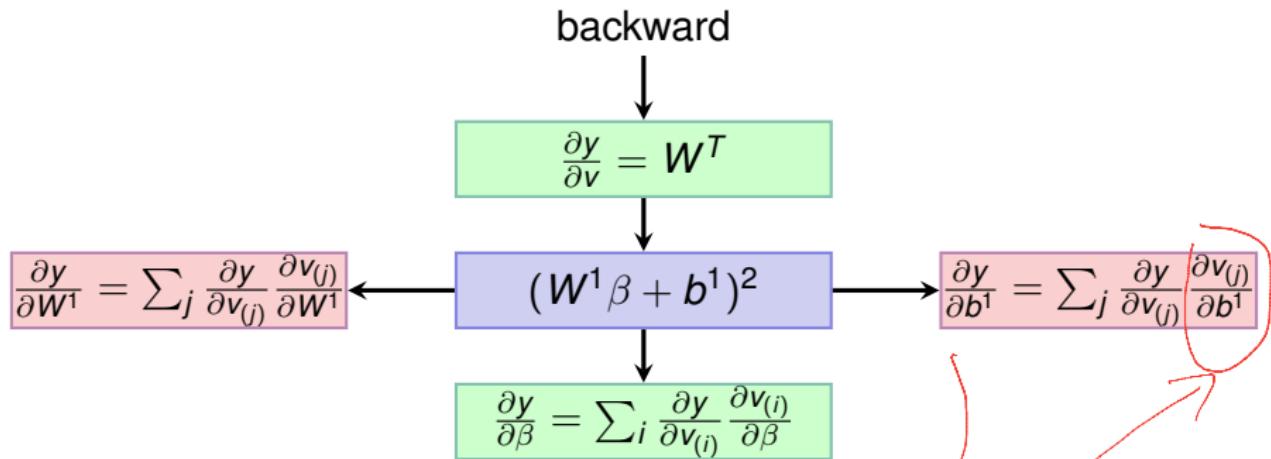
$$z = W^1 \beta + b^1$$

$$v_{(j)} = (e_j^T W^1 \beta + b_{(j)}^1)^2 \stackrel{\text{chain rule}}{=} (\underline{e_j^T W^1 \beta} + \underline{e_j^T b^1})^2$$

$$\Rightarrow \frac{\partial v_{(j)}}{\partial b^1} = 2 \underline{e_j^T W^1 \beta} \cdot e_j$$

$$\frac{\partial v_{(j)}}{\partial z_j} = \frac{\partial (z_j)^2}{\partial z_j} \cdot \frac{\partial z_j}{\partial b^1}$$

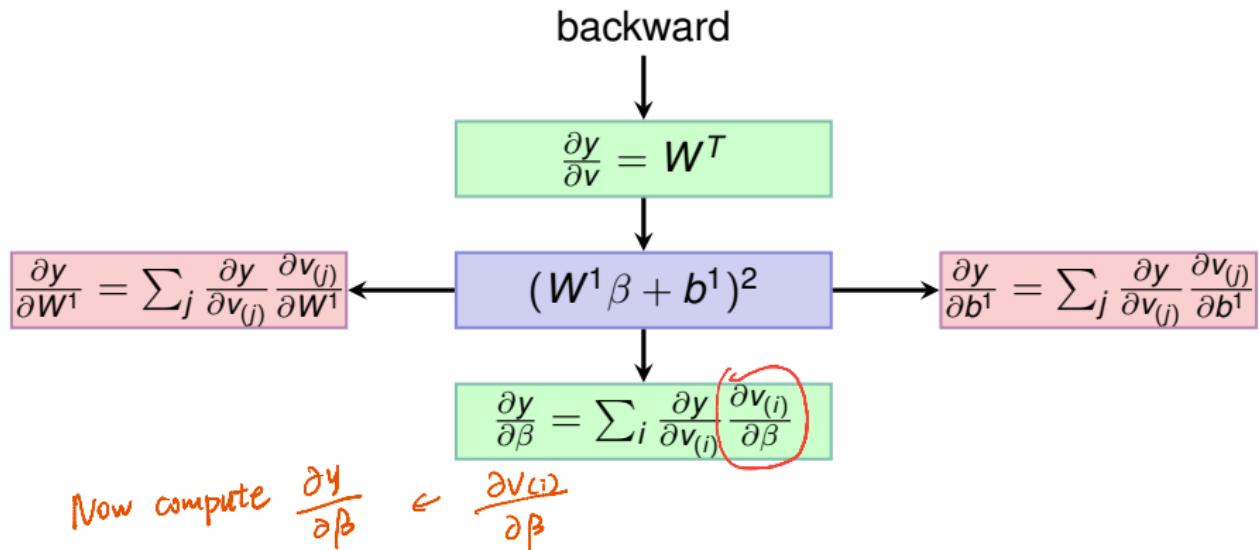
$$\frac{\partial v_{(j)}}{\partial b^1} = 2z_{(j)} e_j$$



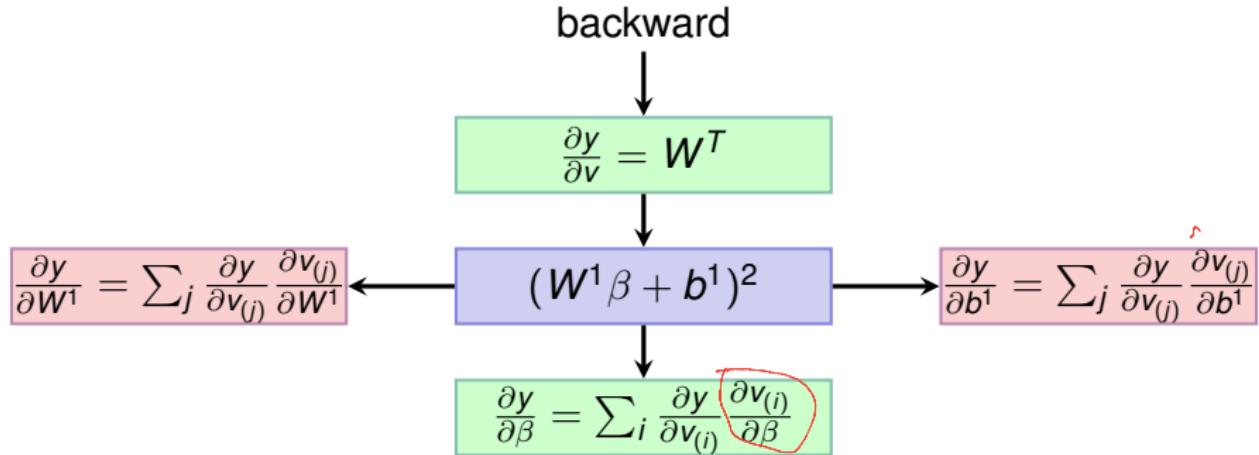
$$z = W^1 \beta + b^1$$

plug in $\frac{\partial v_{(j)}}{\partial b^1} = 2z_{(j)} e_j$

$$\begin{aligned} \frac{\partial y}{\partial b^1} &= \sum_j 2W_{(j)} z_{(j)} e_{(j)} \\ &= 2Wz \quad W^\top \circ z \end{aligned}$$



$$v_{(j)} = (\mathbf{e}_j^T W^1 \beta + b_{(j)}^1)^2$$



$$z = W^1 \beta + b^1$$

$$v_{(j)} = (e_j^T W^1 \beta + b_{(j)}^1)^2$$

$$= \frac{\partial (z_{(j)})^2}{\partial z_{(j)}} \frac{\partial z_{(j)}}{\partial \beta}$$

$$= 2z_{(j)} (e_j^T W^1)^T$$

$$\frac{\partial v_{(j)}}{\partial \beta} = 2z_{(j)} (W^1)^T e_j$$

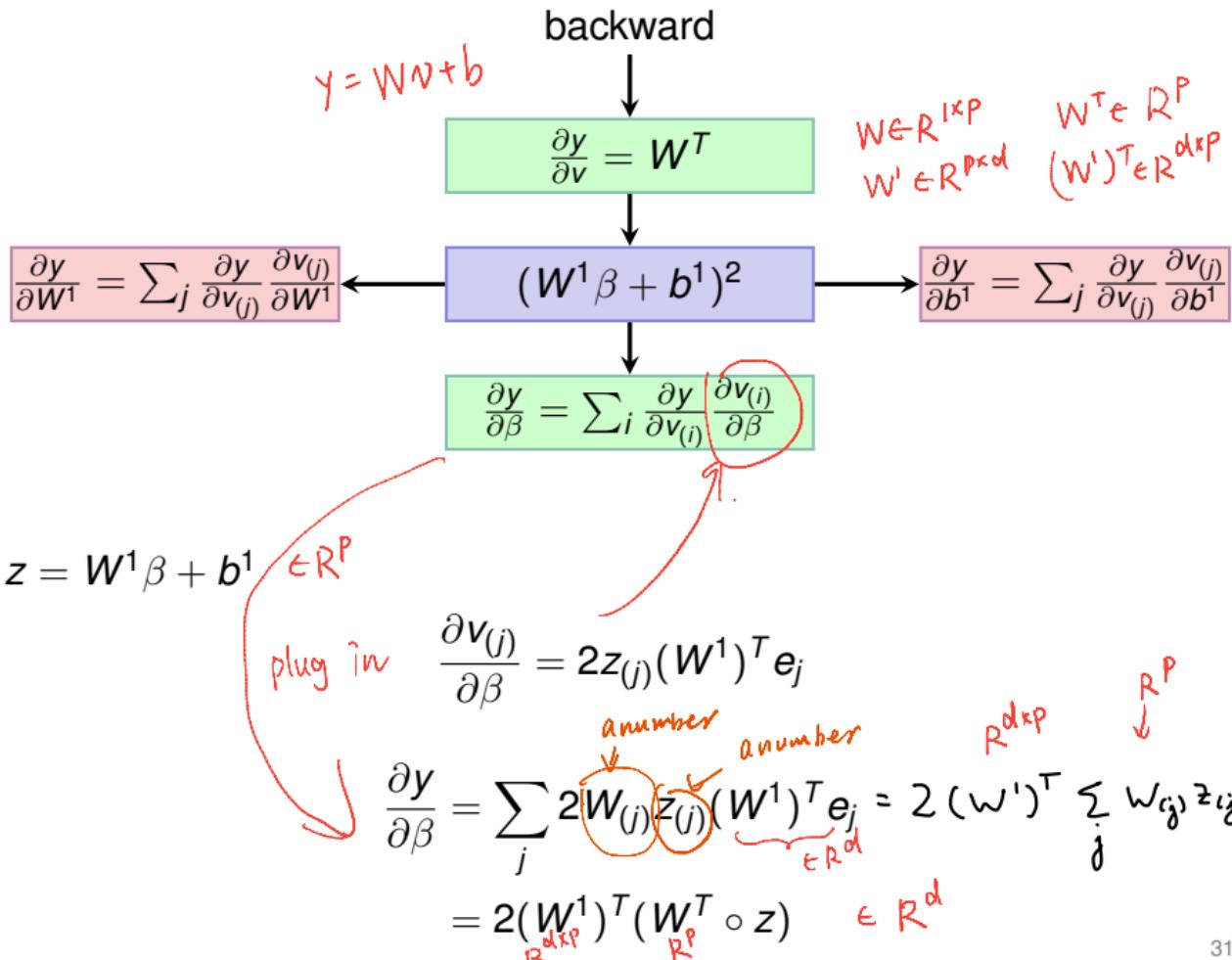
$$e_j \in \mathbb{R}^p \quad W \in \mathbb{R}^{p \times d}$$

$$e_j^T \in \mathbb{R}^{1 \times p}$$

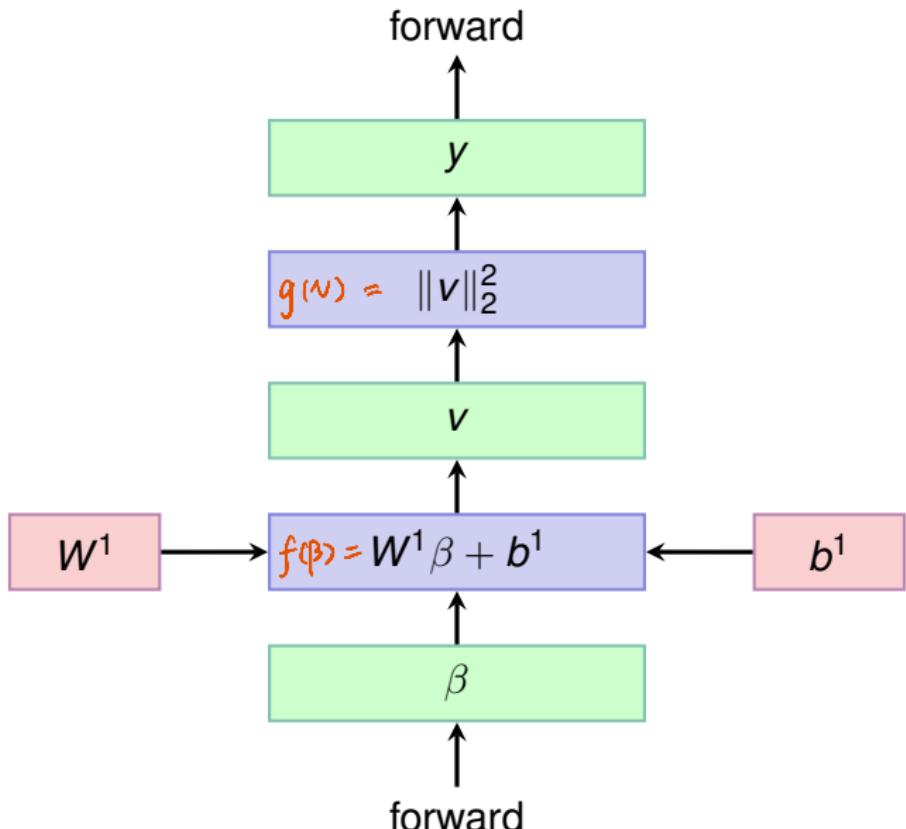
$$(e_j^T W) \in \mathbb{R}^d$$

$$\frac{\partial z_{(j)}}{\partial \beta} \in \mathbb{R}^d$$

等价的维度要和 β 一样



Another example



Mathematically

$$W^1 \in \mathbb{R}^{p \times d}$$

$$b^1 \in \mathbb{R}^p$$

identical to linear regression $y = \|W^1\beta + b^1\|_2^2$

Mathematically

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

$$W^1 \in \mathbb{R}^{p \times d}$$

$$b^1 \in \mathbb{R}^p \quad v \in \mathbb{R}^p$$

$$y = \|W^1\beta + b^1\|_2^2 \quad \epsilon \in \mathbb{R}$$

Already know proof next slide

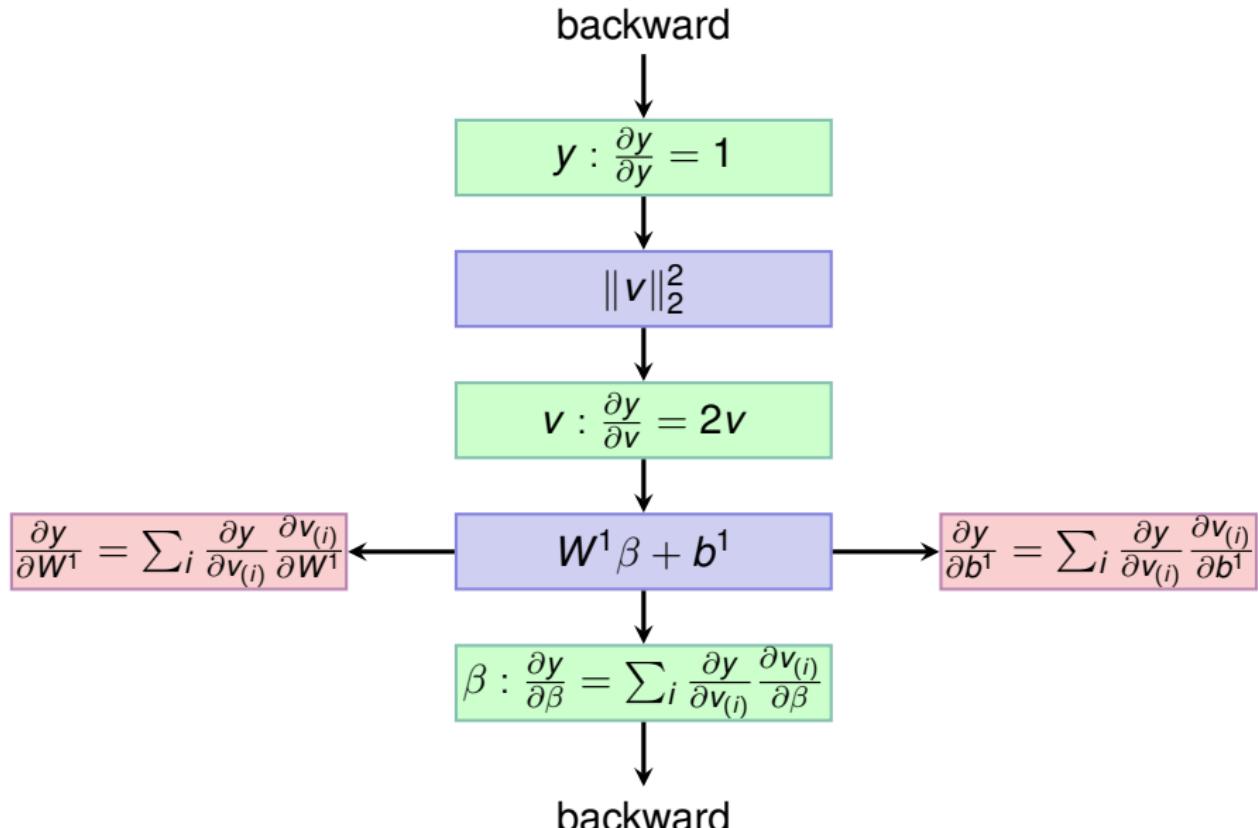
$$\frac{\partial y}{\partial \beta} = 2(W^1)^T(W^1\beta + b^1)$$

$$\frac{\partial y}{\partial b^1} = 2I^T(W^1\beta + b^1) \quad I: \text{identity matrix} \quad y = \|W^1\beta + Ib^1\|_2^2$$

$$\frac{\partial y}{\partial b^1} = 2(W^1\beta + b^1)$$

$$\frac{\partial y}{\partial W^1} = 2(W^1\beta + b^1)\beta^T \quad \text{maybe not this one}$$

Graphically



Computing gradients

$$v = W^1 \beta + b^1$$

$$\frac{\partial y}{\partial v} = 2v$$

加 e_i , v 是步子梯度的
维度和自变量相同

i-th row of W^1

$$\frac{\partial v_{(i)}}{\partial \beta} = W_{(i,:)}^1 = (W^1)^T e_i$$

$$\frac{\partial v_{(i)}}{\partial b^1} = ? e_i$$

$$\frac{\partial v_{(i)}}{\partial W^1} = e_i \beta^T$$

if $AB = BA$ 可交换
then A commute with B

plug in



$$\frac{\partial y}{\partial \beta} = \sum_i 2v_{(i)} W_{(i,:)}^1 = (W^1)^T v$$

$$= \sum_i 2v_{(i)} W^{1T} e_i = 2W^{1T} \sum_i v_i e_i$$

W^1 independent of i
can commute

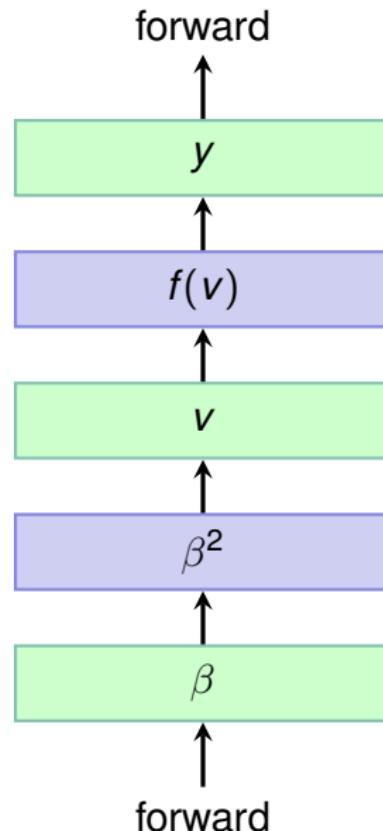
$$\frac{\partial y}{\partial b^1} = \sum_i 2v_{(i)} e_i = 2v$$

$$\frac{\partial y}{\partial W^1} = \sum_i 2v_{(i)} e_i \beta^T = 2v \beta^T$$

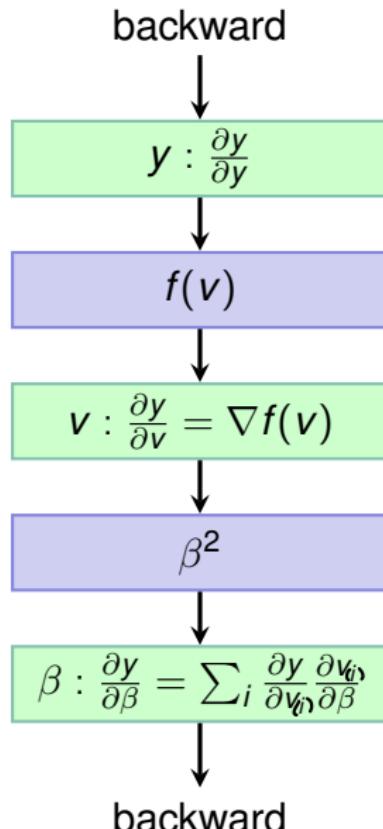
Final example *more easy*

$y = f(\beta^2)$ where $f : \mathbb{R}^d \mapsto \mathbb{R}$

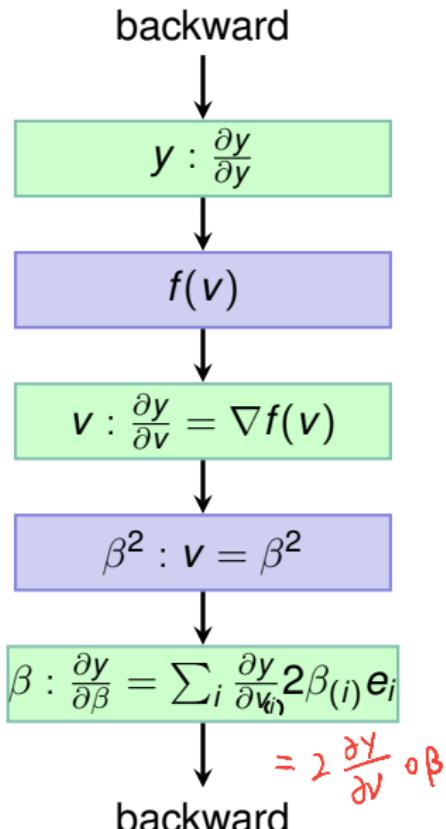
Final example



Final example



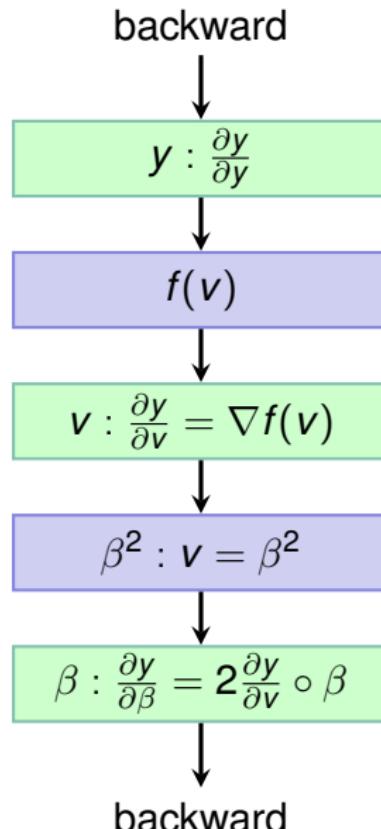
Final example



$$\frac{\partial V_i}{\partial \beta} = 2 \beta_{(i)} e_i$$

$$= 2 \frac{\partial y}{\partial v} \circ \beta$$

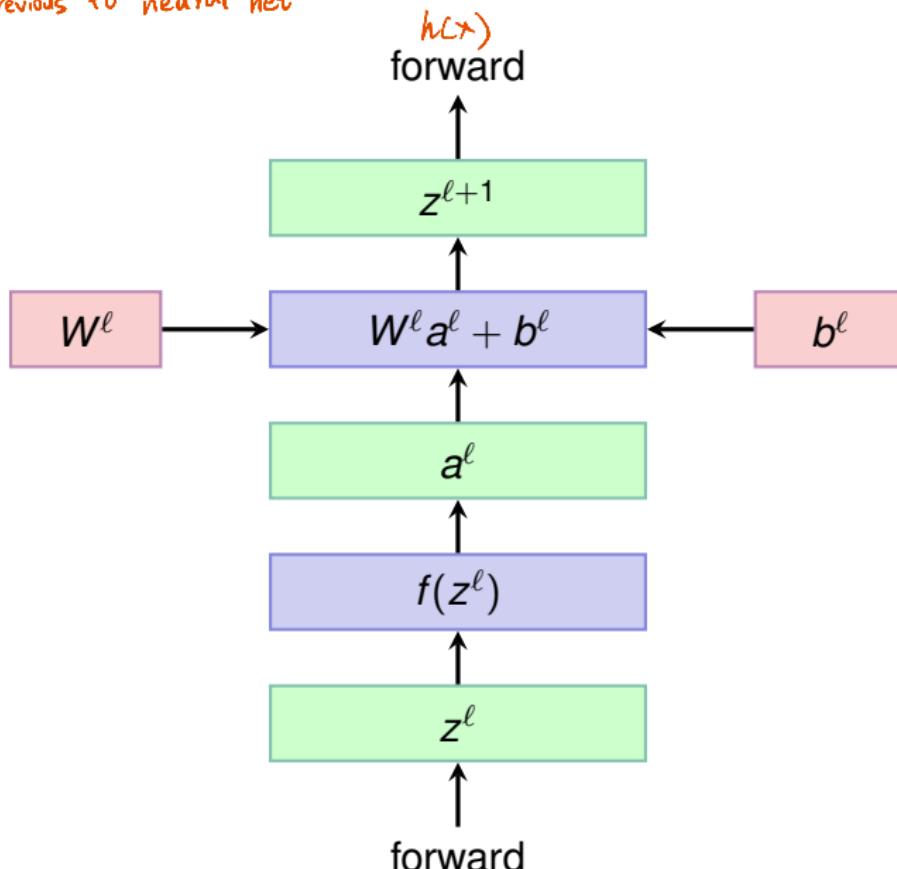
Final example



Neural network: Forward pass

(info)

Now apply previous to neural net.



Backpropagation

求3个梯度
of computing gradient with respect to different params

$h(x)$ backward

$$z_{(i)}^{(l+1)} = e_i^T (W^l a^l + b^l) \quad \in R^P$$

$$\frac{\partial h(x)}{\partial W^l} = \delta^{l+1} (a^l)^T$$

partial
 \downarrow
 \times
 \downarrow
 P^T

$$\delta^{l+1} = \frac{\partial h(x)}{\partial z^{l+1}}$$

a number

$$W^l a^l + b^l$$

$$\frac{\partial h(x)}{\partial b^l} = \delta^{l+1}$$

$$\begin{aligned} \frac{\partial h}{\partial W^l} &= \sum_{(i)} \frac{\partial h}{\partial z_{(i)}^{(l+1)}} \frac{\partial z_{(i)}^{(l+1)}}{\partial W^l} \\ &= \sum_{(i)} \delta_{(i)}^{(l+1)} e_i a^l {}^T \end{aligned}$$

$$= \left(\sum_{(i)} \delta_{(i)}^{(l+1)} e_i \right) a^l {}^T$$

$$= \delta^{(l+1)} [a^{(l)}]^T$$

$$\frac{\partial h(x)}{\partial a^l} = (W^l)^T \delta^{l+1}$$

$$f(z^l)$$

$$\delta^l = \frac{\partial h(x)}{\partial z^l} = f'(z^l) \circ [(W^l)^T \delta^{l+1}]$$

backward

$$\begin{aligned} \frac{\partial h}{\partial a^l} &= \sum_{(i)} \frac{\partial h}{\partial z_{(i)}^{(l+1)}} \frac{\partial z_{(i)}^{(l+1)}}{\partial a^l} \\ &= \sum_{(i)} \delta_{(i)}^{(l+1)} (W^l)^T e_i \end{aligned}$$

$$\begin{aligned} &= (W^l)^T \sum_{(i)} \delta_{(i)}^{(l+1)} e_i \\ &= \sum_{(i)} \delta_{(i)}^{(l+1)} [a^{(l)}]^T \end{aligned}$$

$$\begin{aligned} \frac{\partial h}{\partial z^l} &= \sum_j \frac{\partial h}{\partial a_{(j)}^l} \frac{\partial a_{(j)}^l}{\partial z^l} \\ &= \sum_j \frac{\partial h}{\partial a_{(j)}^l} \delta'((z^l)_{(j)}) e_j \end{aligned}$$

$$= \sum_j \frac{\partial h}{\partial a_{(j)}^l} \delta'((z^l)_{(j)}) e_j$$

Conclusion

Gradient methods will behave differently based on geometry of problem

Conclusion

Gradient methods will behave differently based on geometry of problem

Neural networks: powerful methods to approximate non-linear functions

Backpropagation: computationally efficient method to compute gradients of NN

↓
gradient descent