

Back-propagation

CPT_S 434/534 Neural network design and application

Determining model parameters

- Computational complexity for the analytical solution?

$$\nabla_w f(w) = \frac{1}{n} \sum_{i=1}^n x_i' w x_i - y_i x_i \rightarrow 0 \Rightarrow XX'w^* - XY = 0 \Rightarrow w^* = (XX')^{-1}XY$$

- Matrix multiplication:

$$XX': d \times n \times d \quad XY: d \times n \quad (XX')^{-1}XY: d \times d \times n \rightarrow O(d^2n)$$

- Inverse of a matrix:

$$(XX')^{-1}: O(d^{2.373})$$

- Total complexity

$$O(d^2n + d^{2.373})$$

Optimization for machine learning

- **First-order** algorithms (commonly used and researched in machine learning)
 - Gradient descent
 - Momentum methods
 - Stochastic variants
 - Hessian vector products
 -

Optimization for machine learning

- **First-order** algorithms (commonly used and researched in machine learning)
 - **Gradient** descent
 - Momentum methods
 - Stochastic variants
 - Hessian vector products
 -

First-order → need to compute **gradients**

Optimization for machine learning

- **First-order** algorithms (commonly used and researched in machine learning)
 - **Gradient** descent
 - Momentum methods
 - Stochastic variants
 - Hessian vector products
 -

$$h_S^{\text{ERM}} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_S(h).$$

First-order → need to compute **gradients**

Optimization for machine learning

- **First-order** algorithms (commonly used and researched in machine learning)
 - **Gradient** descent
 - Momentum methods
 - Stochastic variants
 - Hessian vector products
 -

$$h_S^{\text{ERM}} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_S(h).$$

$$\min_{w_1, \dots, w_K} \frac{1}{n} \sum_{i=1}^n \left(\underbrace{- \sum_{k=1}^K y_{ik} \cdot \log(f(w_k; x_i))}_{\triangleq L_i(W) \text{ (see Section 2)}} \right),$$

First-order → need to compute **gradients**

Optimization for machine learning

- **First-order** algorithms (commonly used and researched in machine learning)
 - **Gradient** descent
 - Momentum methods
 - Stochastic variants
 - Hessian vector products
 -

First-order → need to compute **gradients**

$$h_S^{\text{ERM}} = \operatorname{argmin}_{h \in \mathcal{H}} \widehat{R}_S(h).$$

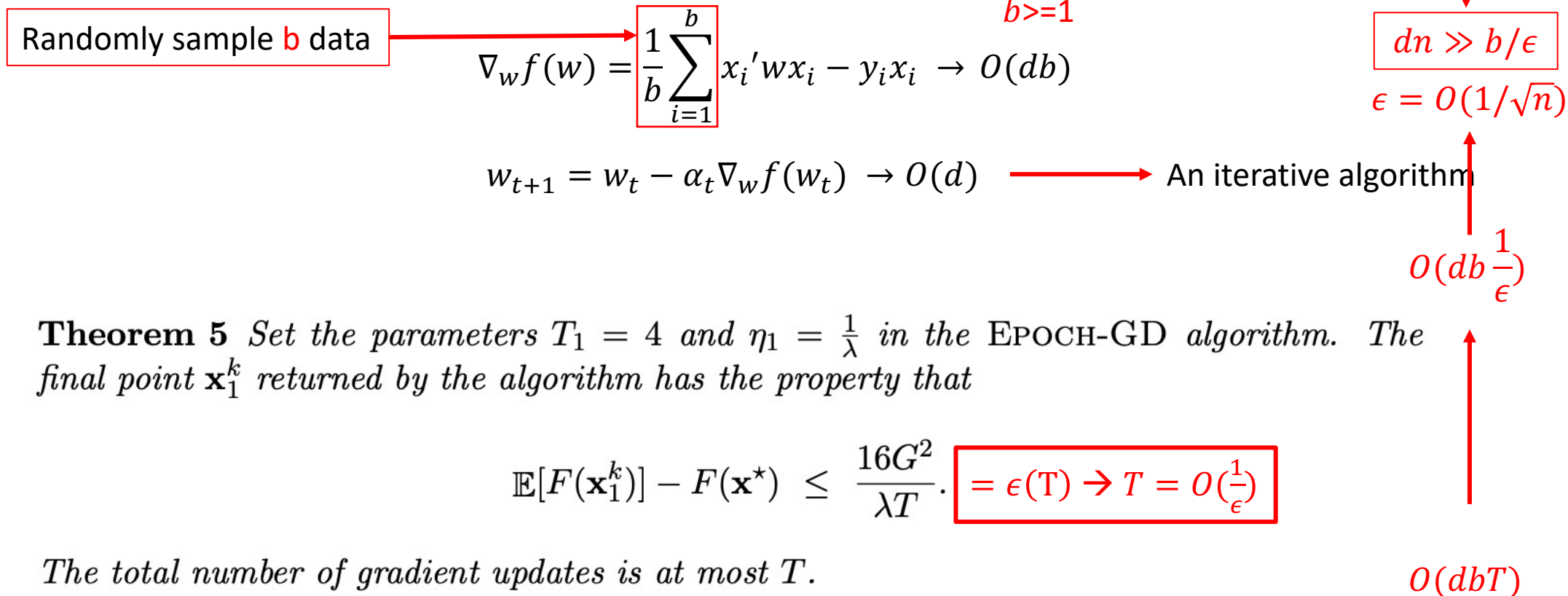
$$\min_{w_1, \dots, w_K} \frac{1}{n} \sum_{i=1}^n \left(\underbrace{- \sum_{k=1}^K y_{ik} \cdot \log(f(w_k; x_i))}_{\triangleq L_i(W) \text{ (see Section 2)}} \right),$$

$$\frac{\partial F}{\partial w_k} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial w_k} + \lambda w_k.$$

Gradients for updating

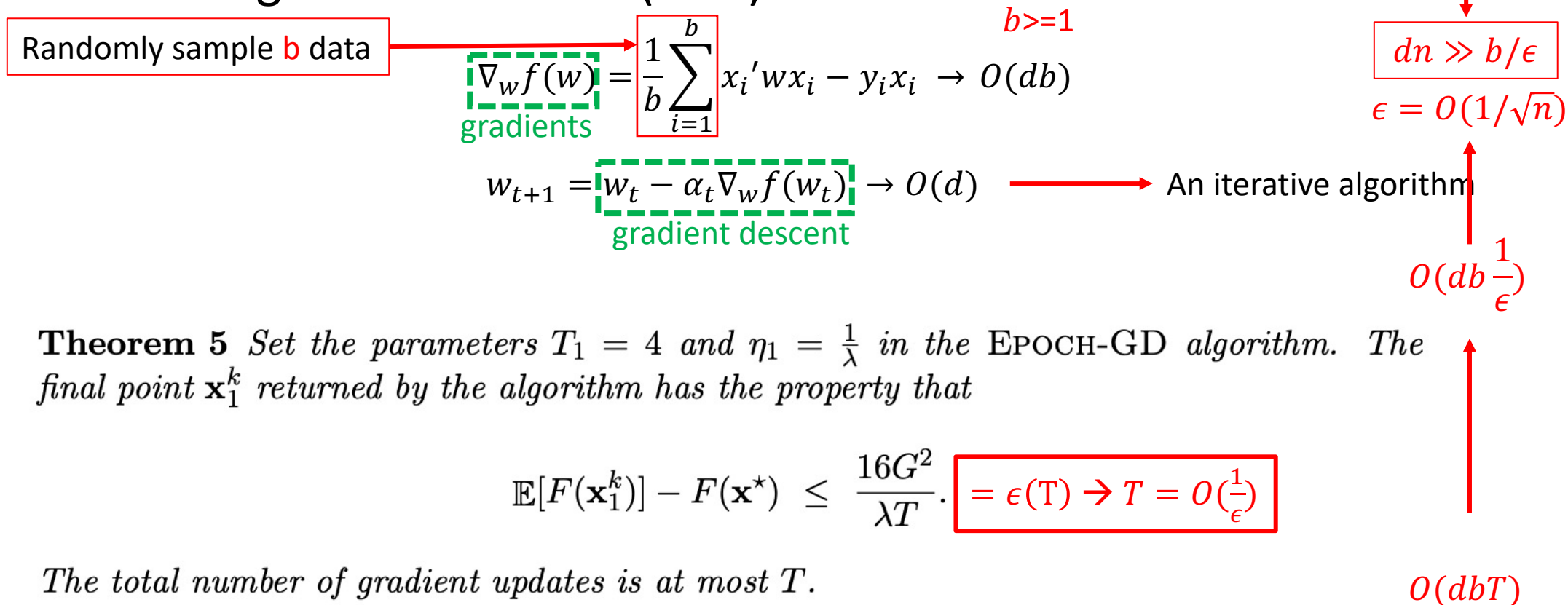
Determining model parameters

- Stochastic gradient descent (SGD)



Determining model parameters

- Stochastic gradient descent (SGD)



How to compute gradient?

$$f(x) \rightarrow ?$$

How to compute gradient?

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

How to compute gradient?

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

What if we have composition structure?

f and *g* both have their own parameters

x is the parameter of function *g*

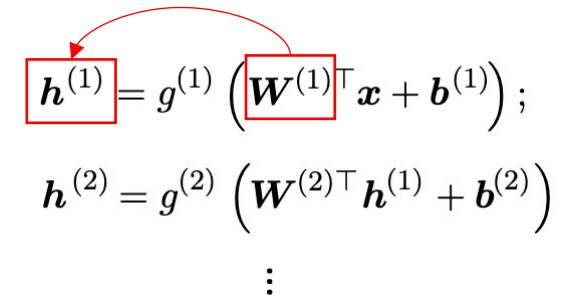
$$f(g(x)) \rightarrow ?$$

How to compute gradient?

What if we have composition structure?
f and *g* both have their own parameters
x is the parameter of function *g*

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$


$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

How to compute gradient?

What if we have composition structure?
f and *g* both have their own parameters
x is the parameter of function *g*

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \boxed{h^{(1)}} &= g^{(1)} \left(\boxed{W^{(1)}}^\top \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \boxed{h^{(2)}} &= g^{(2)} \left(\boxed{W^{(2)}}^\top \boxed{h^{(1)}} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$

$$\frac{dz}{dx} = \boxed{\frac{dz}{dy}} \boxed{\frac{dy}{dx}}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \boxed{h^{(1)}} &= g^{(1)} \left(\boxed{W^{(1)}}^\top \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \boxed{h^{(2)}} &= g^{(2)} \left(\boxed{W^{(2)}}^\top \boxed{h^{(1)}} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$
 $= f(g(x))$

$$\frac{dz}{dx} = \boxed{\frac{dz}{dy}} \boxed{\frac{dy}{dx}}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$
 $= f(g(x))$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$f(y) = \log(y)$$

$$g(x) = 1 + e^{-x}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$
 $= f(g(x))$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$f(y) = \log(y)$$

$$\nabla f(y) = 1/y$$

$$g(x) = 1 + e^{-x}$$

$$\nabla g(x) = -e^{-x}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$
 $= f(g(x))$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\begin{aligned} f(y) &= \log(y) & g(x) &= 1 + e^{-x} \\ \nabla f(y) &= 1/y & \nabla g(x) &= -e^{-x} \end{aligned}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$\boxed{y = g(x)} \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \boxed{h^{(1)}} &= g^{(1)} \left(\boxed{W^{(1)}}^\top \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \boxed{h^{(2)}} &= g^{(2)} \left(\boxed{W^{(2)}}^\top \boxed{h^{(1)}} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$
 $= f(g(x))$

$$\frac{dz}{dx} = \boxed{\frac{dz}{dy}} \boxed{\frac{dy}{dx}}$$

$$f(y) = \log(y)$$

$$\nabla f(y) = 1/y$$

$$g(x) = 1 + e^{-x}$$

$$\nabla g(x) = -e^{-x}$$

$$\nabla h(x) = \frac{-e^{-x}}{1 + e^{-x}}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$\boxed{y = g(x)} \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \boxed{h^{(1)}} &= g^{(1)} \left(\boxed{W^{(1)}}^\top x + b^{(1)} \right); \\ \boxed{h^{(2)}} &= g^{(2)} \left(\boxed{W^{(2)}}^\top \boxed{h^{(1)}} + b^{(2)} \right) \\ &\vdots \end{aligned}$$

An example:

$$\begin{aligned} h(x) &= \log(1 + e^{-x}) \\ &= f(g(x)) \end{aligned}$$

$$\frac{dz}{dx} = \boxed{\frac{dz}{dy} \frac{dy}{dx}}$$

$$f(y) = \log(y)$$

$$\nabla f(y) = 1/y$$

$$g(x) = 1 + e^{-x}$$

$$\nabla g(x) = -e^{-x}$$

$$\nabla h(x) = \frac{-e^{-x}}{1 + e^{-x}}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

- Chain rule of calculus

$$\boxed{y = g(x)} \text{ and } z = f(g(x)) = f(y)$$

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \boxed{h^{(1)}} &= g^{(1)} \left(\boxed{W^{(1)}}^\top \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \boxed{h^{(2)}} &= g^{(2)} \left(\boxed{W^{(2)}}^\top \boxed{h^{(1)}} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

An example: $h(x) = \log(1 + e^{-x})$
 $= f(g(x))$

$$\frac{dz}{dx} = \boxed{\frac{dz}{dy}} \boxed{\frac{dy}{dx}}$$

$f(y) = \log(y)$
 $\nabla f(y) = 1/y$

$g(x) = 1 + e^{-x}$
 $\nabla g(x) = -e^{-x}$

$\nabla h(x) = \frac{-e^{-x}}{1 + e^{-x}}$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$\begin{aligned} \boxed{h^{(1)}} &= g^{(1)} \left(\boxed{W^{(1)}}^\top \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \boxed{h^{(2)}} &= g^{(2)} \left(\boxed{W^{(2)}}^\top \boxed{h^{(1)}} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(\underbrace{f_1(x)}_{x_1} \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(\underbrace{f_1(x)}_{x_1} \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$f_n \left(\dots \left(f_2(x_1) \right) \right) \rightarrow ?$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \underbrace{\left(f_2(f_1(x)) \right)}_{x_2} \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

How to compute gradient?

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

What if we have composition structure?

f and g both have their own parameters

x is the parameter of function g

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

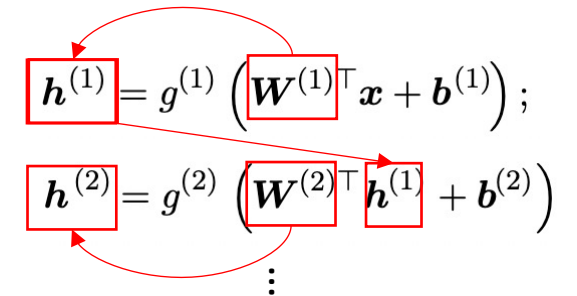
$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

How to compute gradient?

What if we have composition structure?
f and *g* both have their own parameters
x is the parameter of function *g*

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$



$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

Q: is $\frac{dx_n}{dx_1}$ enough to update the model (a lot of layers)?

How to compute gradient?

What if we have composition structure?
f and *g* both have their own parameters
x is the parameter of function *g*

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

Q: is $\frac{dx_n}{dx_1}$ enough to update the model (a lot of layers)?

NO. There are parameters to be determined in each layer

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

Q: is $\frac{dx_n}{dx_1}$ enough to update the model (a lot of layers)?

NO. There are parameters to be determined in each layer

We still need $\frac{dx_n}{dx_i}$, for $i = 1, \dots, n - 1$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

Q: gradient at other hidden layers?

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

How to compute gradient?

What if we have composition structure?
f and *g* both have their own parameters
x is the parameter of function *g*

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

Q: gradient at other hidden layers?

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

$$\frac{dx_n}{dx_i}$$

How to compute gradient?

What if we have composition structure?
 f and g both have their own parameters
 x is the parameter of function g

$$f(x) \rightarrow \nabla f(x) = \frac{df}{dx}$$

$$f(g(x)) \rightarrow ?$$

$$\begin{aligned} \mathbf{h}^{(1)} &= g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right); \\ \mathbf{h}^{(2)} &= g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right) \\ &\vdots \end{aligned}$$

- Chain rule of calculus (generalize to multi-dimensional cases)

$$f_n \left(\dots \left(f_2 \left(f_1(x) \right) \right) \right) \rightarrow ?$$

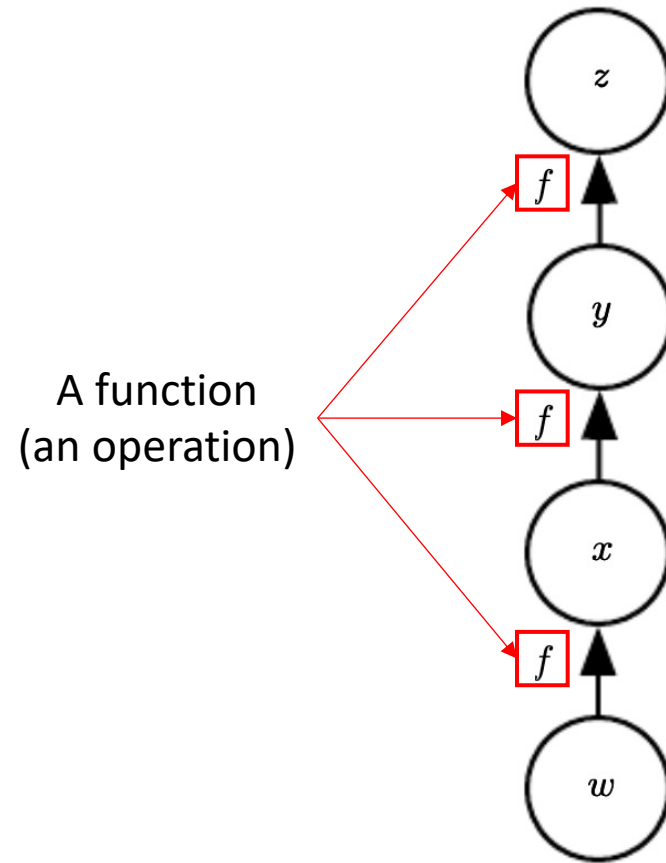
$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

Q: gradient at other hidden layers?

$$\frac{dx_n}{dx_i} = \frac{dx_n}{dx_{n-1}} \cdot \frac{dx_{n-1}}{dx_{n-2}} \cdot \dots \cdot \frac{dx_{i+1}}{dx_i}$$

Computation graphs



Computation graphs

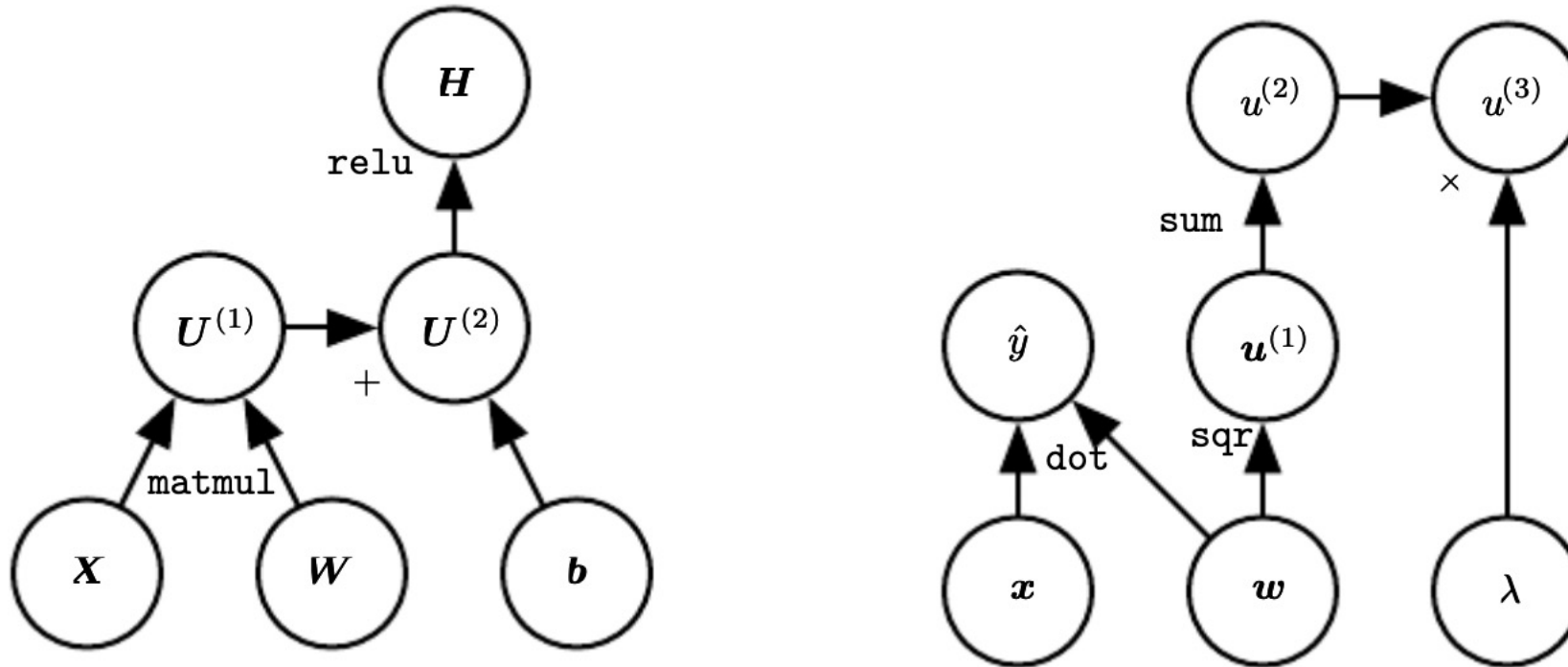


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Computation graphs

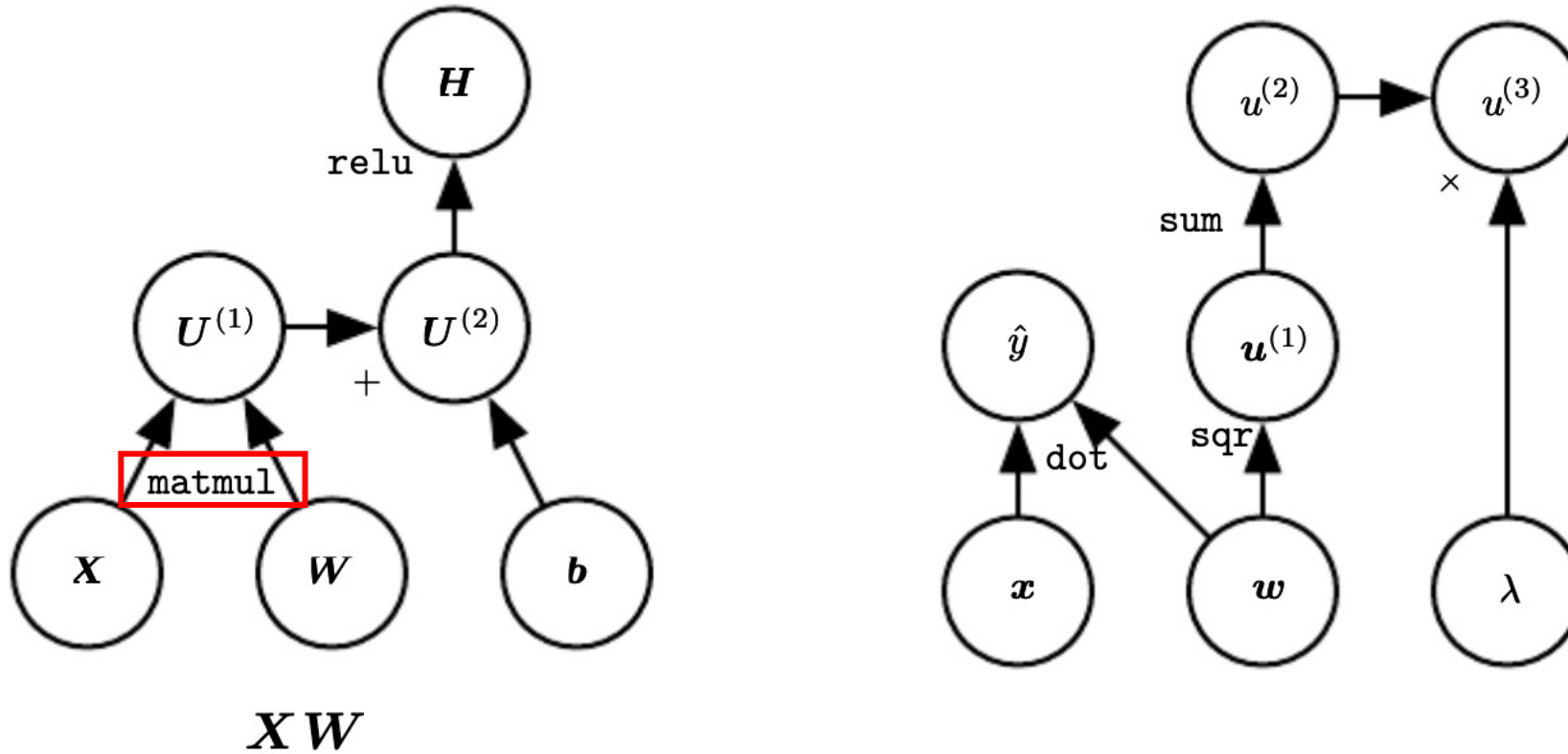


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Computation graphs

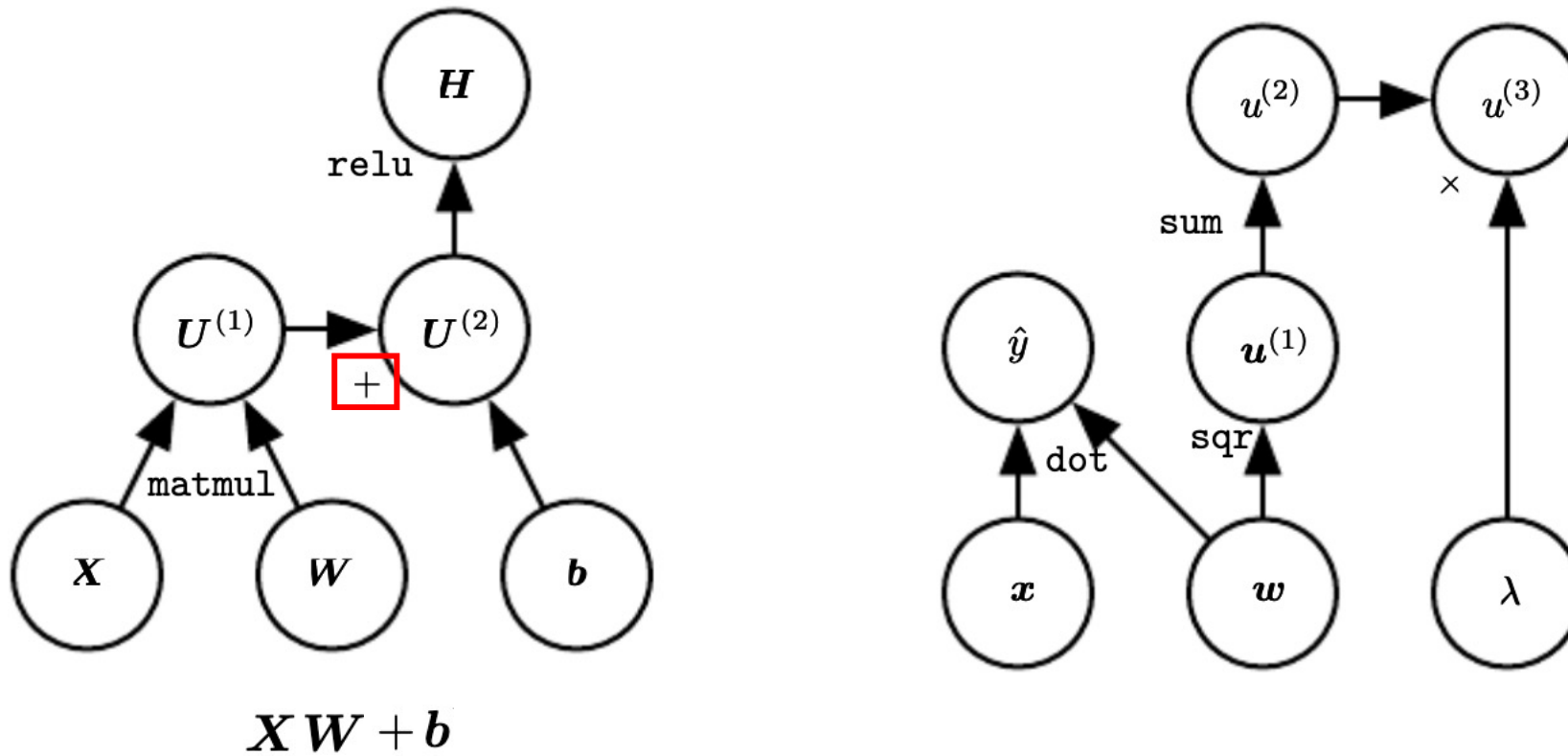


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Computation graphs

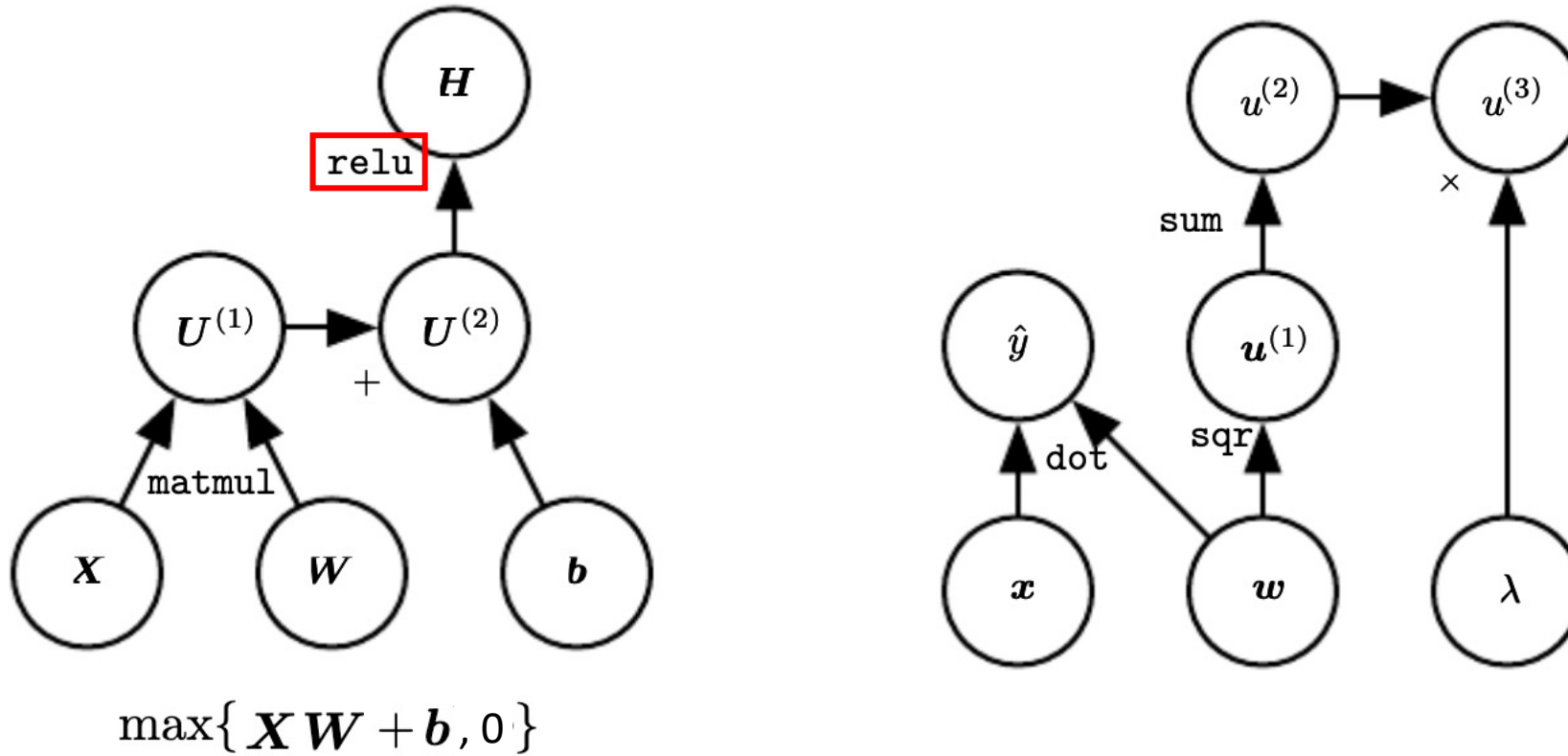


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Computation graphs

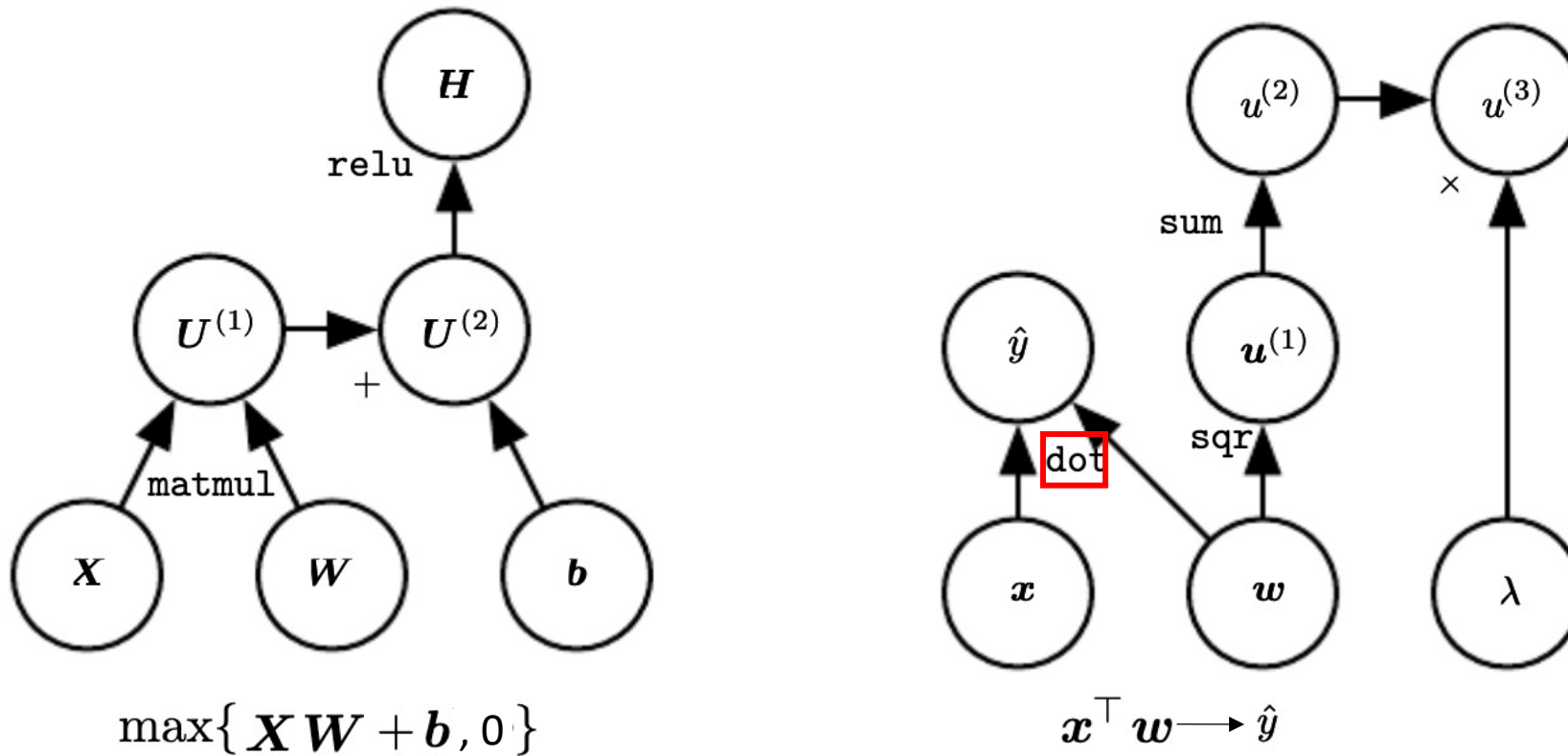


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Computation graphs

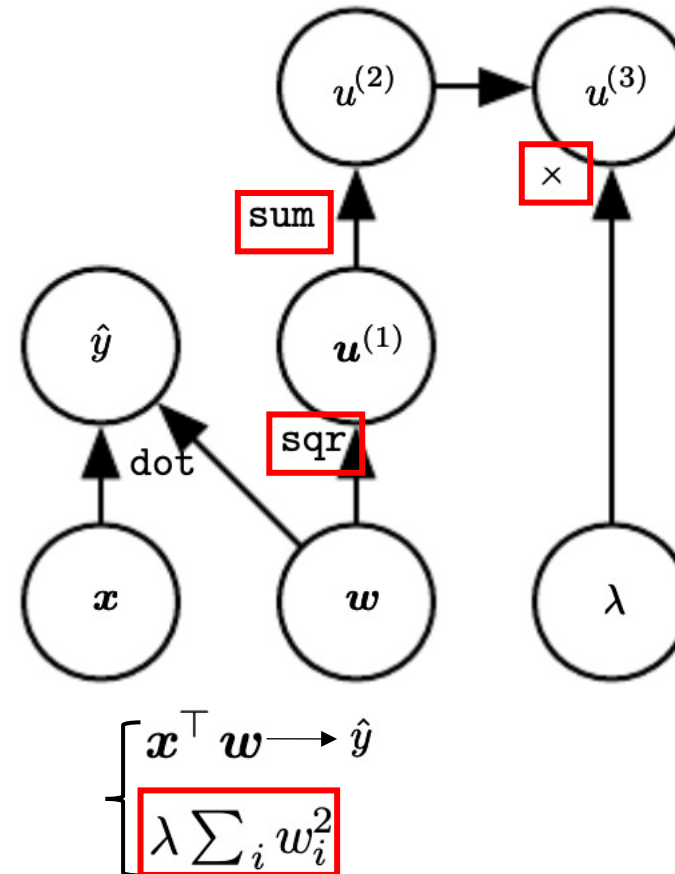
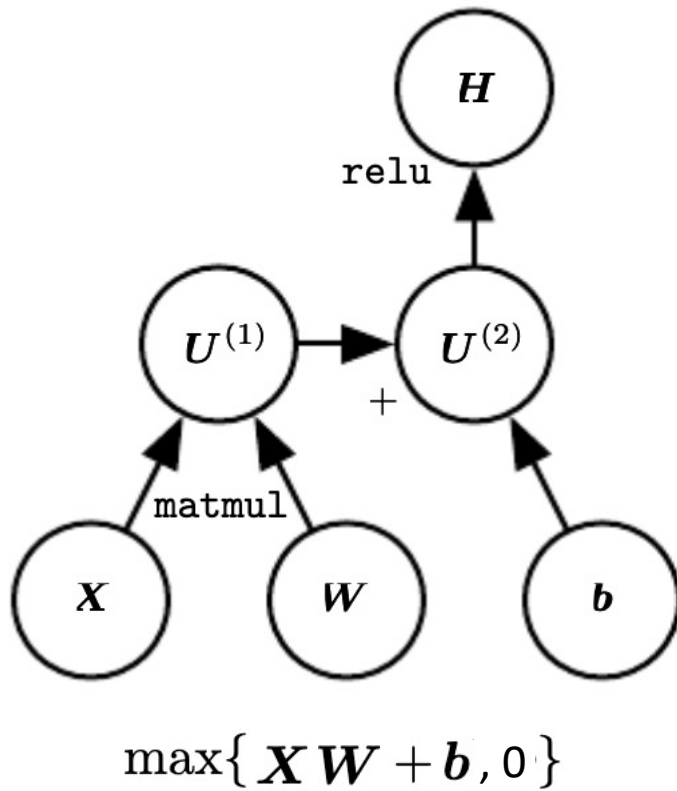


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Computation graphs

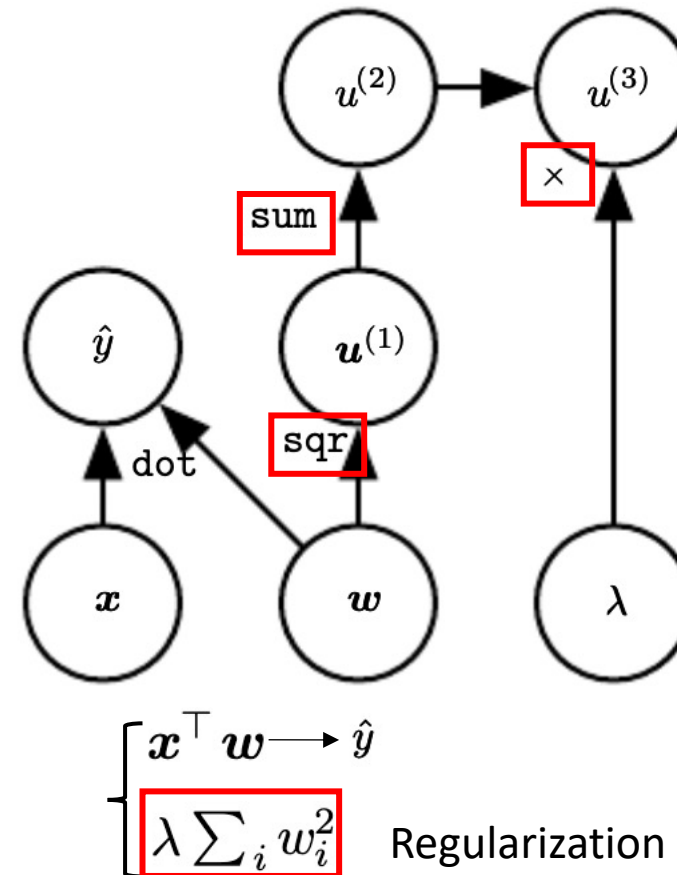
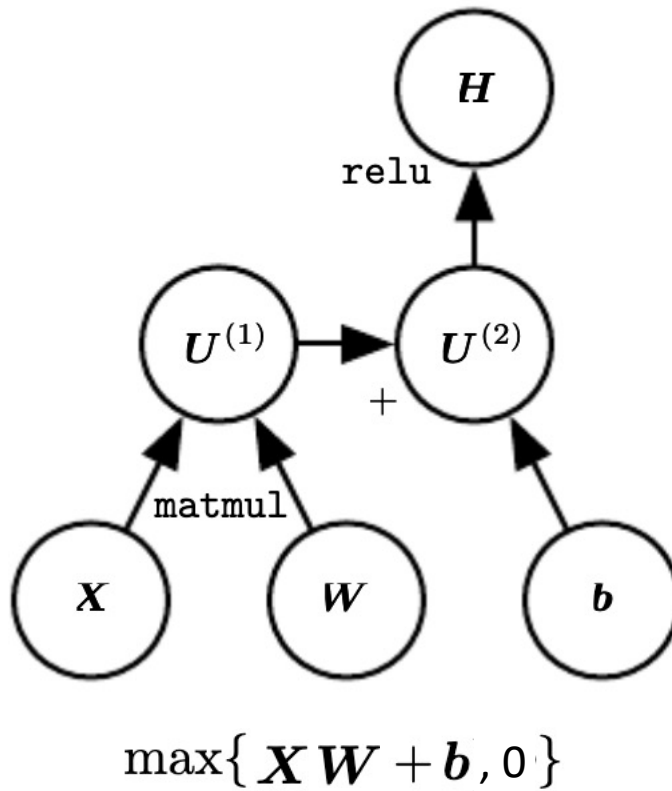


Figure 6.8
“Deep Learning”

It is a precise language to describe structure of operations in neural networks

Forward propagation

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} , the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

$$h^{(l)} \left(\dots h^{(3)} \left(h^{(2)} \left(h^{(1)}(\mathbf{x}) \right) \right) \right)$$

Forward propagation

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

$$h^{(l)} \left(\dots h^{(3)} \left(h^{(2)} \left(h^{(1)} (h^{(0)}) \right) \right) \right)$$

Forward propagation

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

$$h^{(l)} \left(\dots h^{(3)} \left(h^{(2)} \left(h^{(1)} (h^{(0)}) \right) \right) \right)$$

Forward propagation

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

$$h^{(l)} \left(\dots h^{(3)} \left(h^{(2)} \left(h^{(1)} (h^{(0)}) \right) \right) \right)$$

Backward propagation

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y}) \quad \text{Gradient from loss}$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient on the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

Backward propagation

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y}) \quad \text{Gradient from loss}$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient on the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)}) \quad \text{Gradient from activation layer}$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

Backward propagation

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y}) \quad \text{Gradient from loss}$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient on the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)}) \quad \text{Gradient from activation layer}$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta) \quad \text{Gradient from regularization}$$

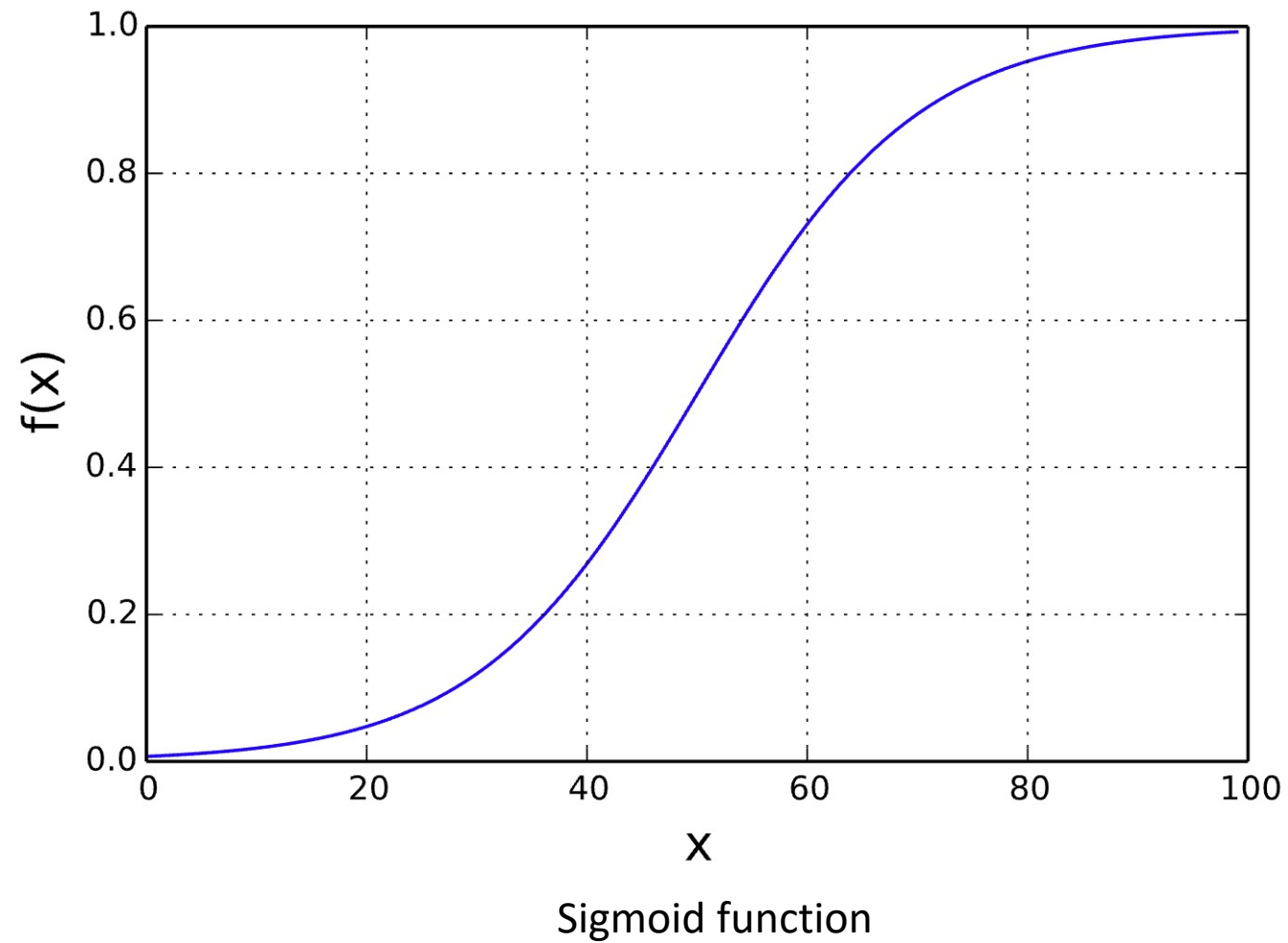
$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta) \quad \text{Gradient from regularization}$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

Gradient vanish



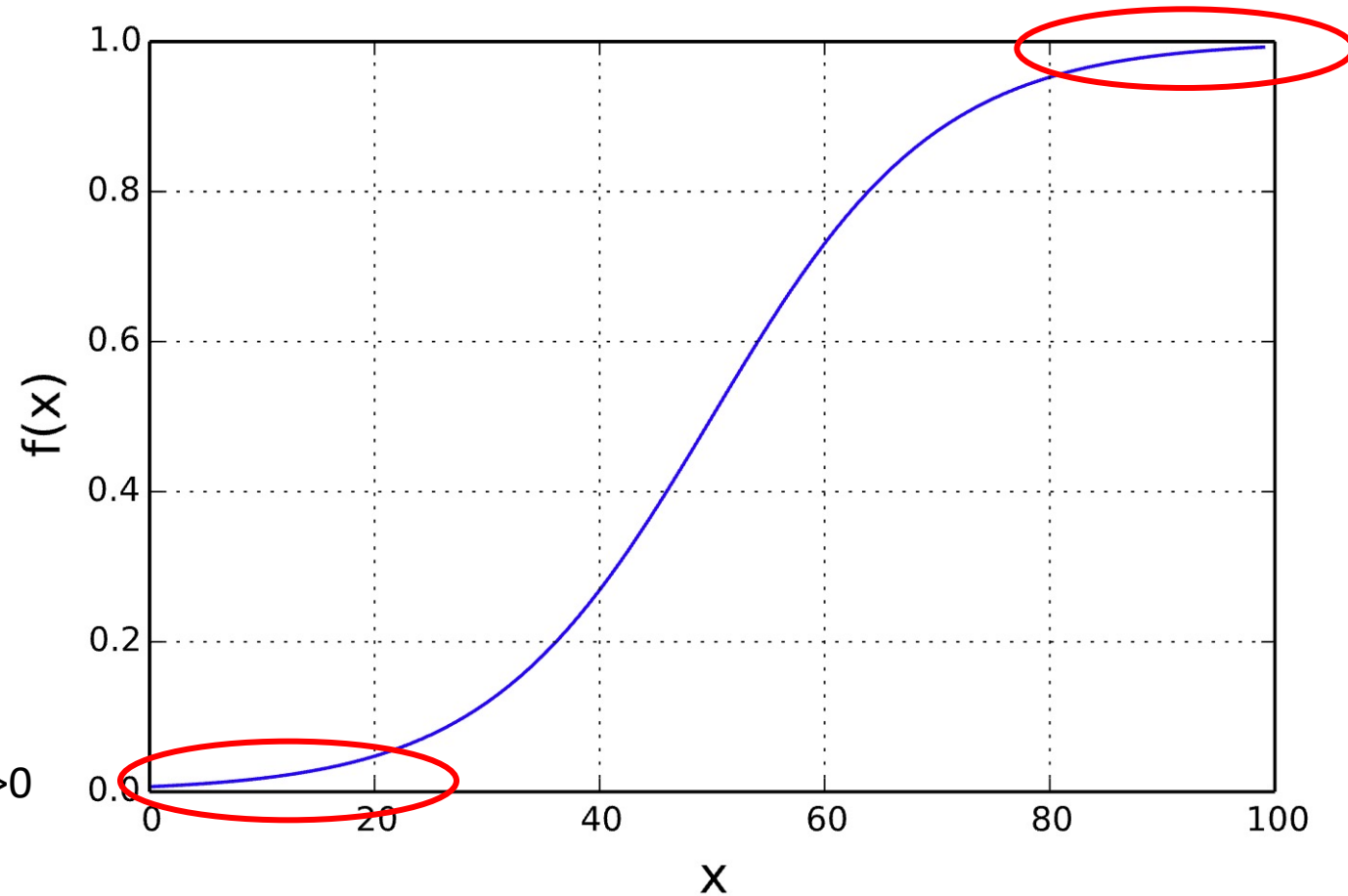
Gradient vanish

$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

gradients $\rightarrow 0$



gradients $\rightarrow 0$

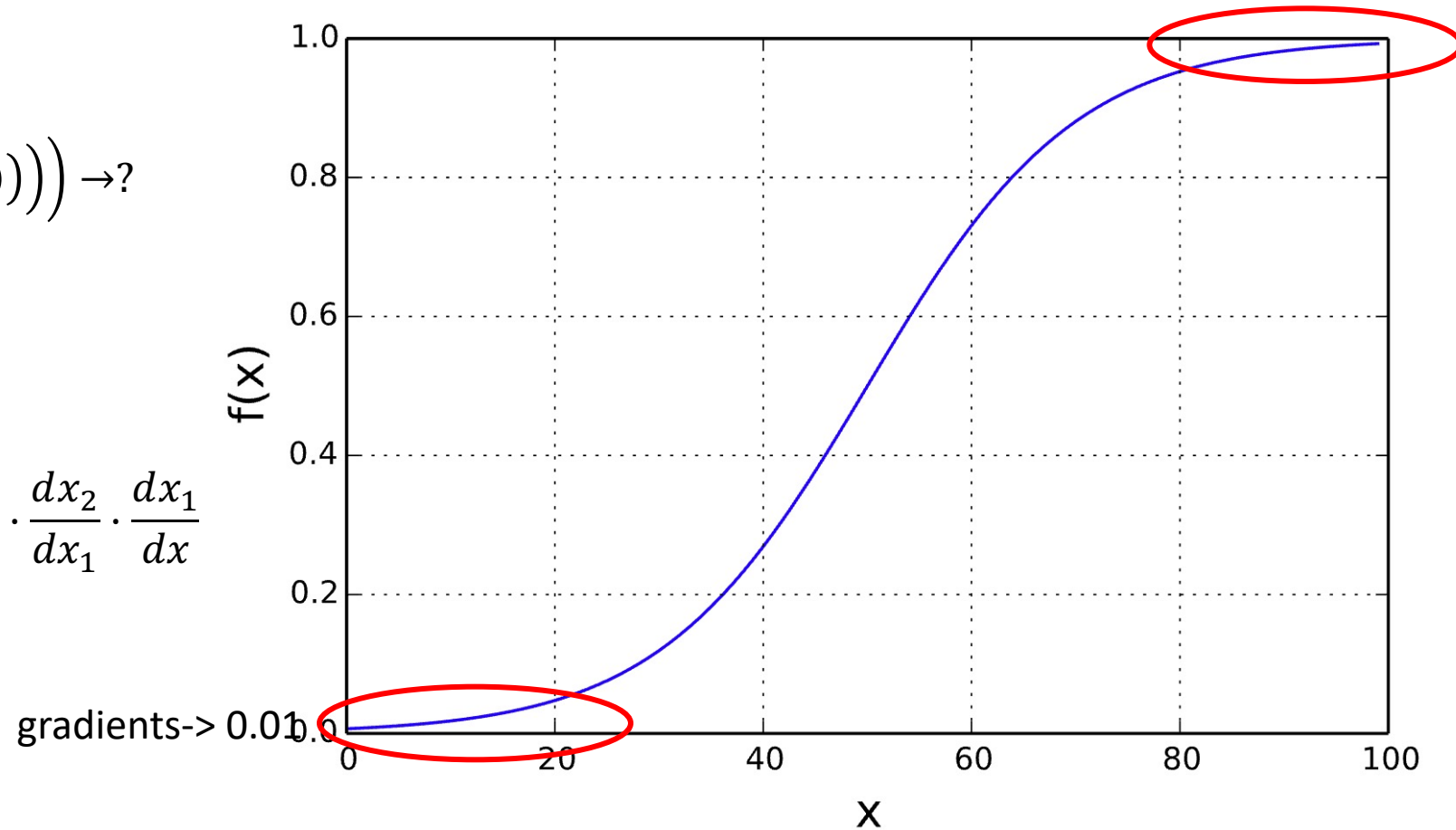
Sigmoid function

Gradient vanish

$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$



gradients $\rightarrow 0.01$

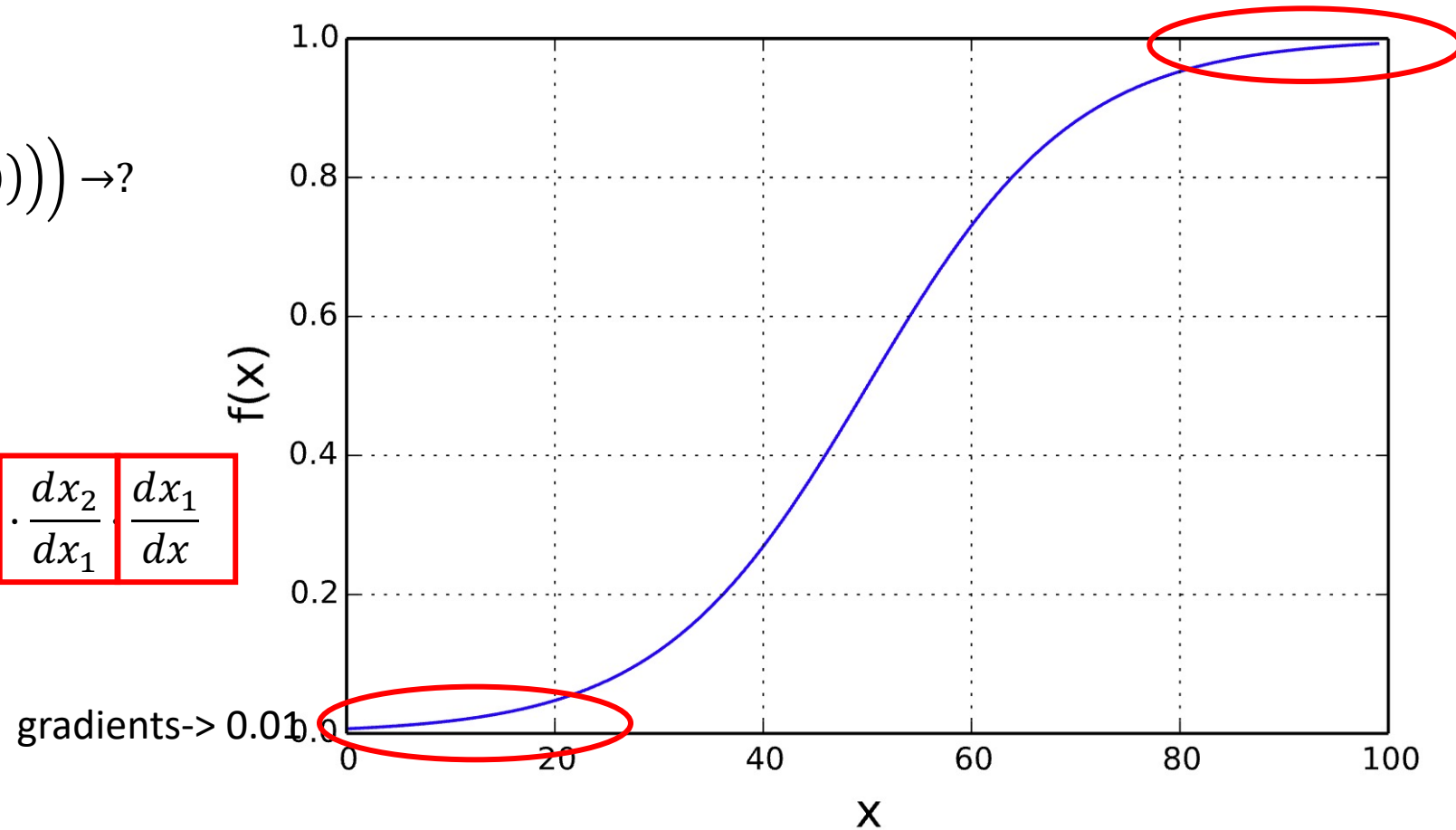
Sigmoid function

Gradient vanish

$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \boxed{\frac{dx_n}{dx_{n-1}}} \cdots \boxed{\frac{dx_2}{dx_1} \frac{dx_1}{dx}}$$



gradients $\rightarrow 0.01$

Sigmoid function

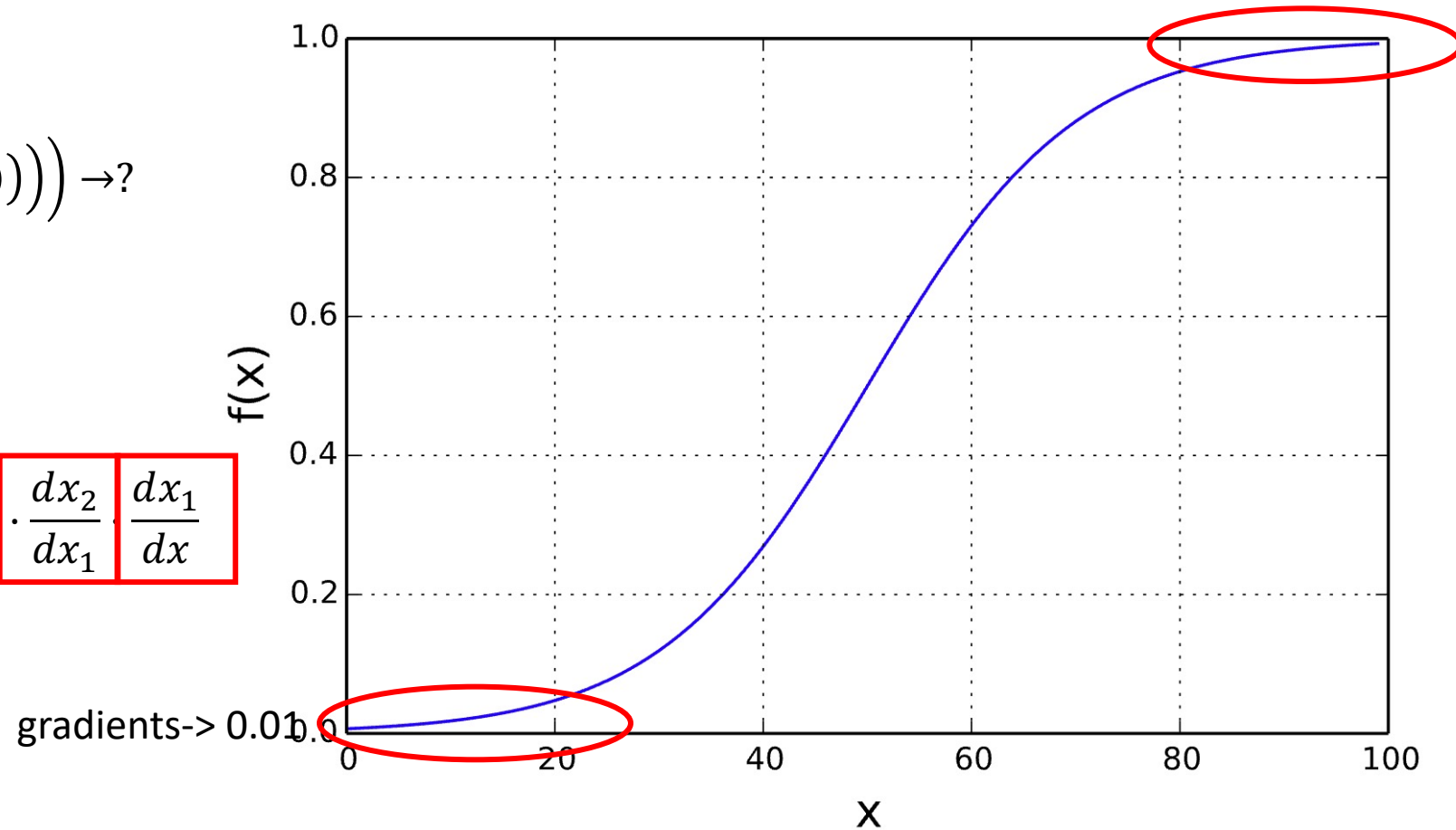
Gradient vanish

$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \boxed{\frac{dx_n}{dx_{n-1}}} \cdots \boxed{\frac{dx_2}{dx_1} \frac{dx_1}{dx}}$$

$\rightarrow 0.01^n$



gradients $\rightarrow 0.01$

gradients $\rightarrow 0.01$

Sigmoid function

Gradient explosion

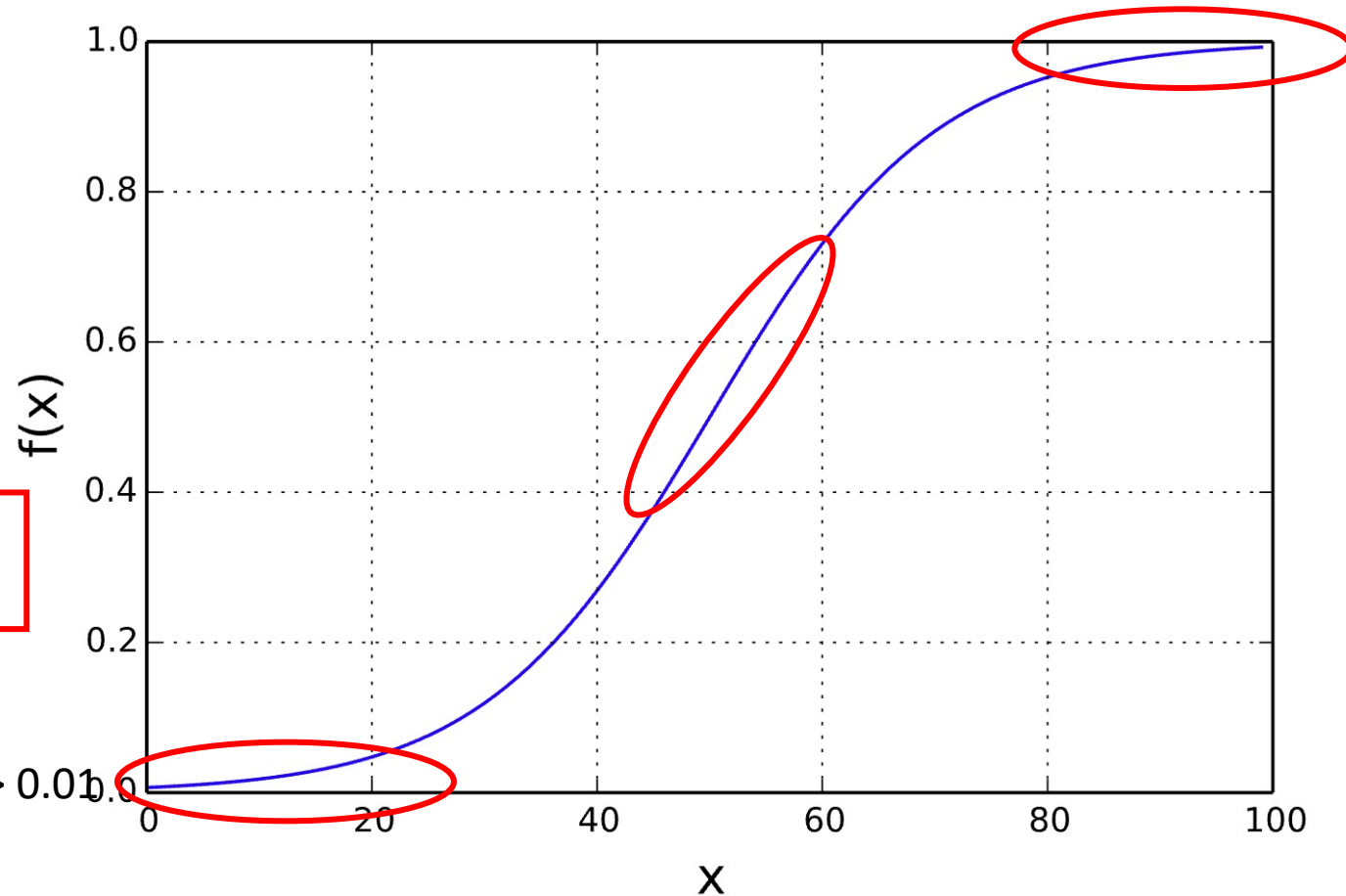
$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \boxed{\frac{dx_n}{dx_{n-1}}} \cdots \boxed{\frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}}$$

$$\rightarrow 1.1^n$$

gradients $\rightarrow 0.01$



gradients $\rightarrow 0.01$

Sigmoid function

Gradient explosion

$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

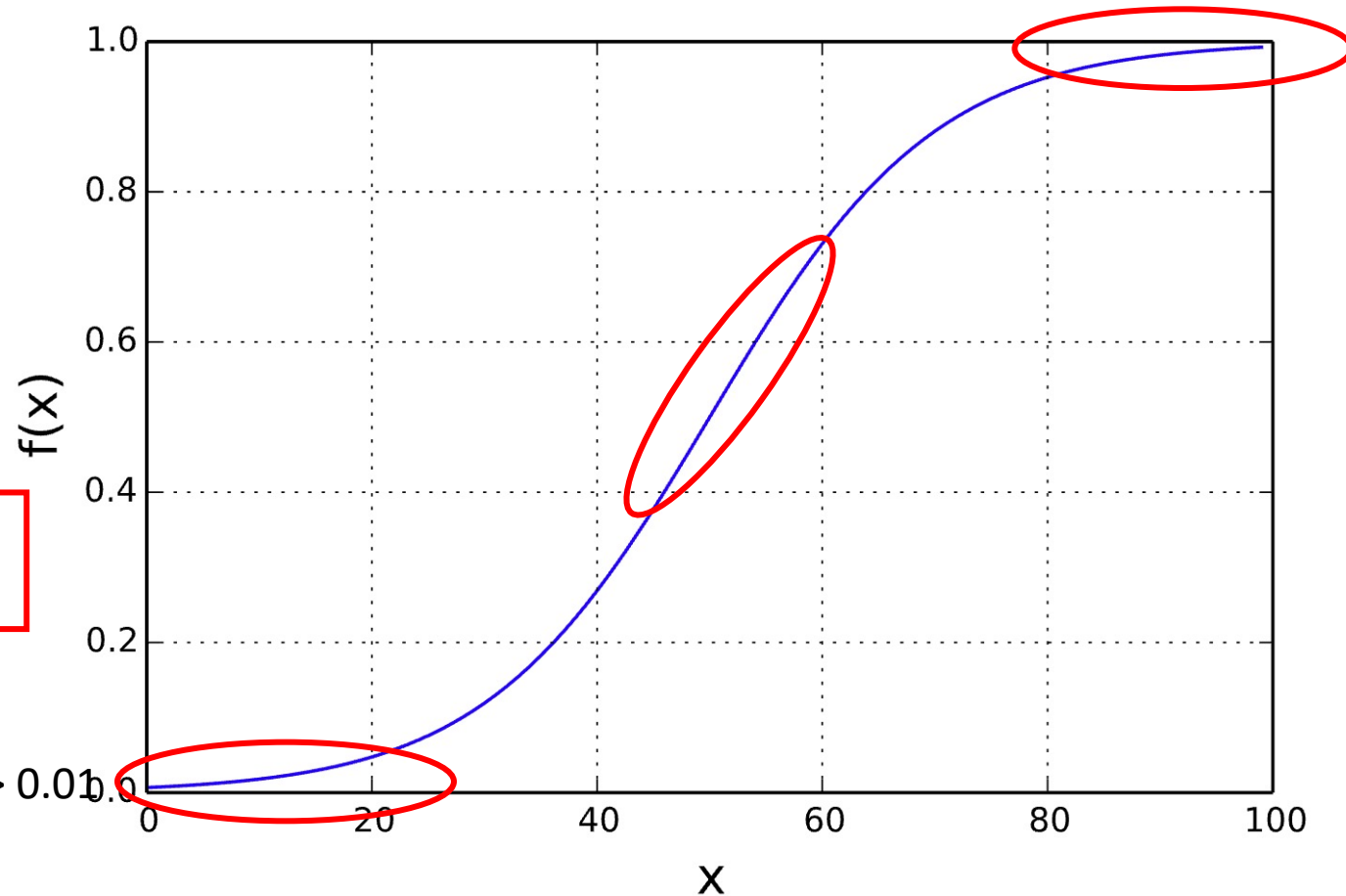
$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx} = \boxed{\frac{dx_n}{dx_{n-1}}} \cdots \boxed{\frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}}$$

$$\rightarrow 1.1^n$$

$$1.1^{100} = 13781$$

gradients $\rightarrow 0.01$



gradients $\rightarrow 0.01$

Sigmoid function