



UNIVERSITÀ  
DI PAVIA

# *Deep Learning*

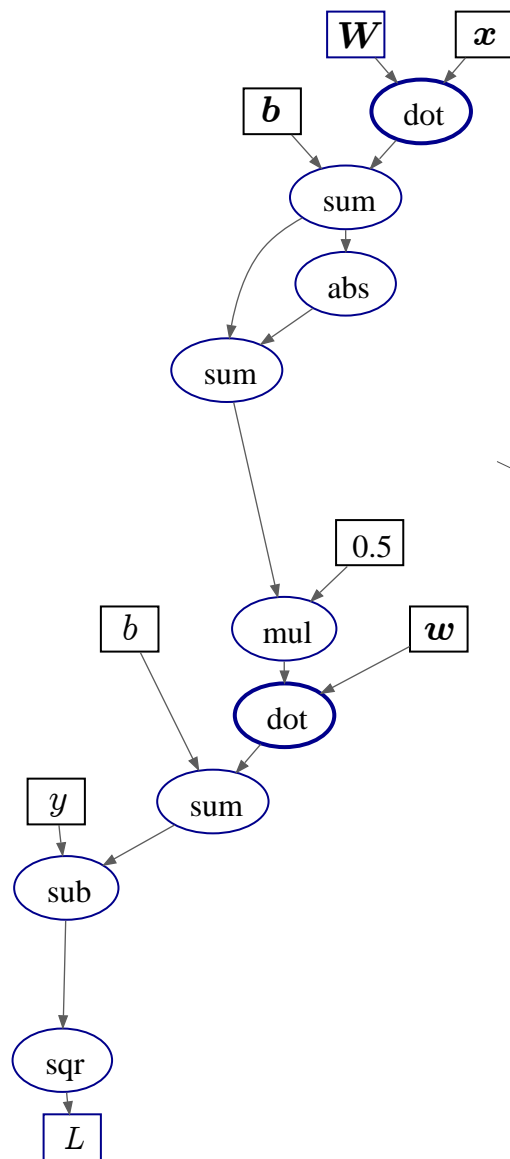
## *03- Flow Graphs & Automatic Differentiation*

Marco Piastra

*This presentation can be downloaded at:*  
<http://vision.unipv.it/DL>

# *Flow Graphs*

# An aside: Flow Graph

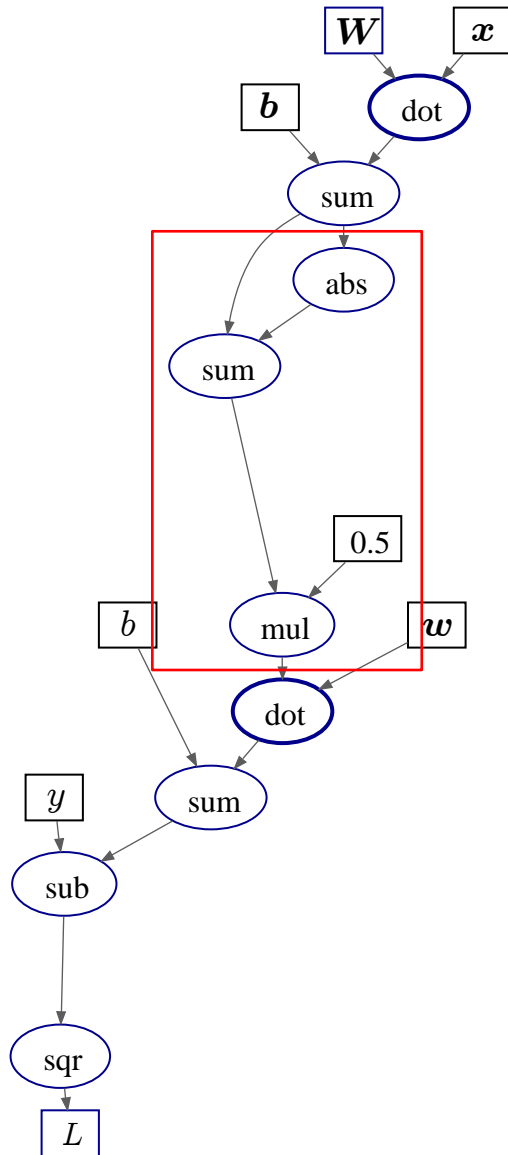


$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Item-wise loss function, FF neural network with ReLU as non-linearity

The above expression translates into this flow graph

# An aside: Flow Graph



$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

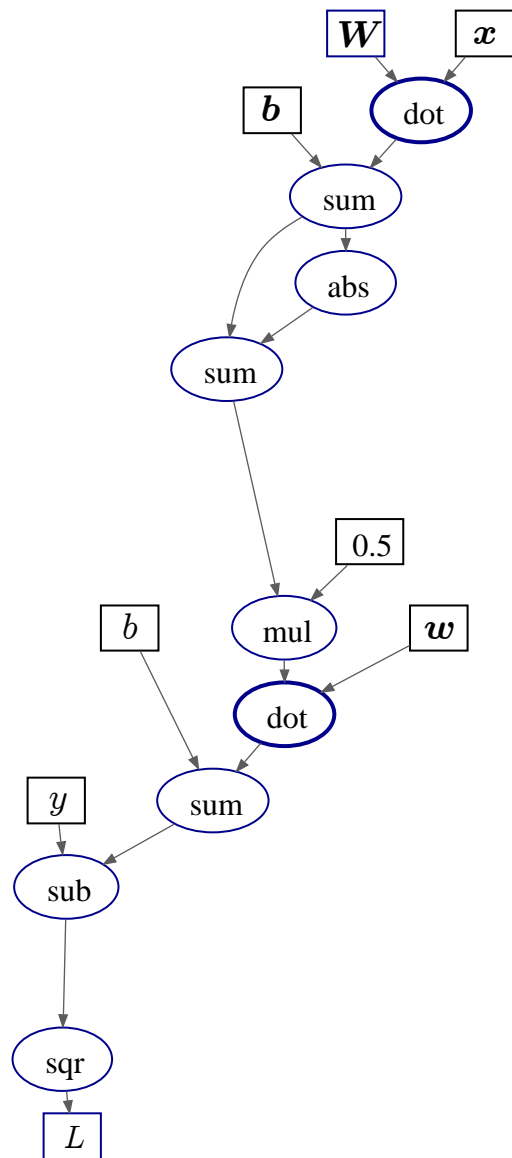
*Item-wise loss function, FF neural network with ReLU as non-linearity*

$$\text{ReLU}(x) := \max(0, x)$$

$$\text{ReLU}(x) = \frac{1}{2}(x + |x|)$$

*(equivalent expression)*

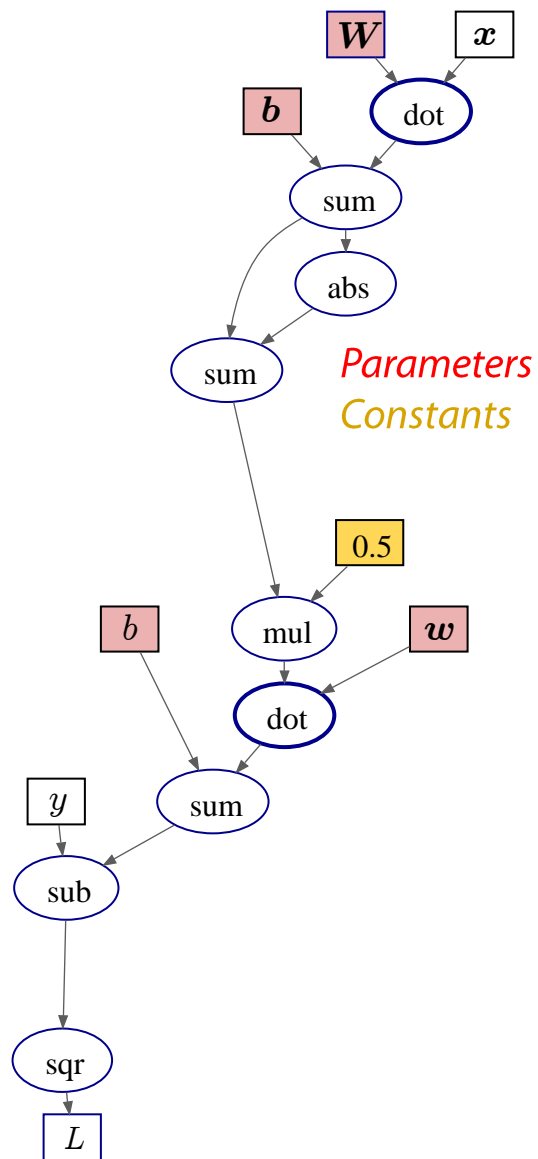
# An aside: Flow Graph



$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

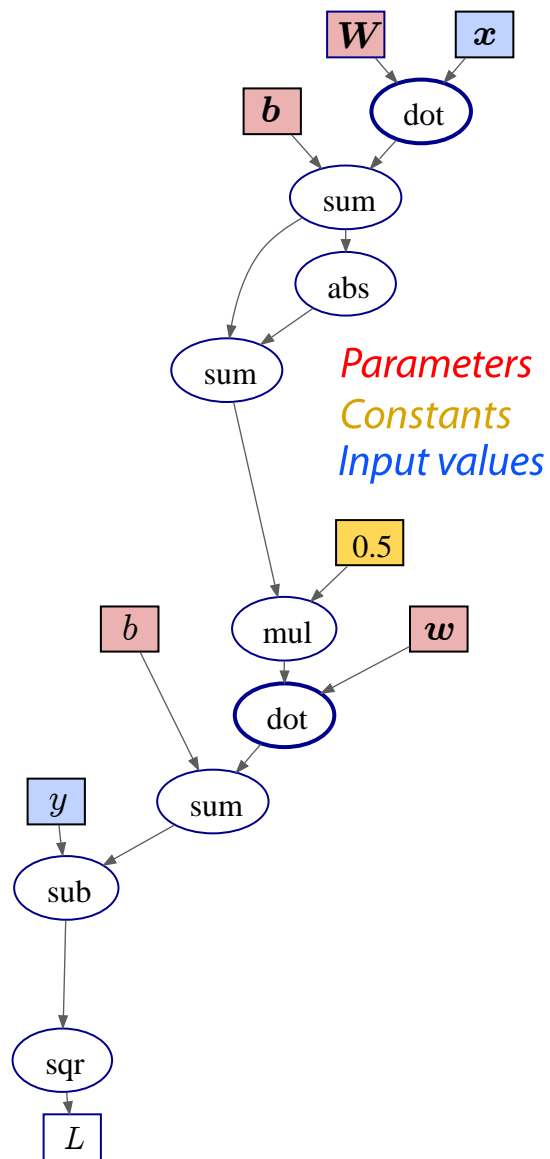
# An aside: Flow Graph



$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Item-wise loss function, FF neural network with ReLU as non-linearity

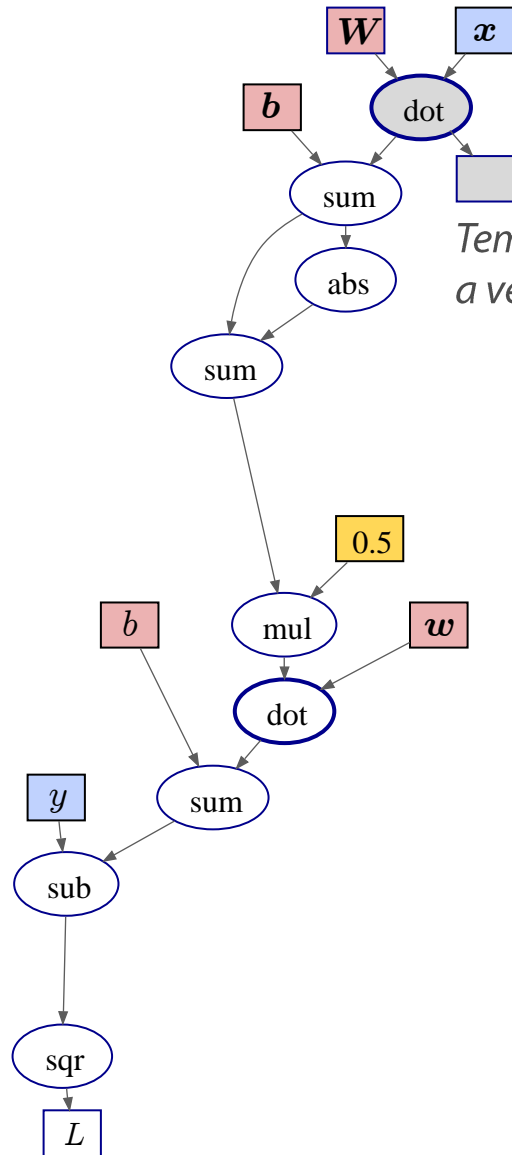
# An aside: Flow Graph



$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Item-wise loss function, FF neural network with ReLU as non-linearity

# An aside: Flow Graph



Temporary value:  
a vector

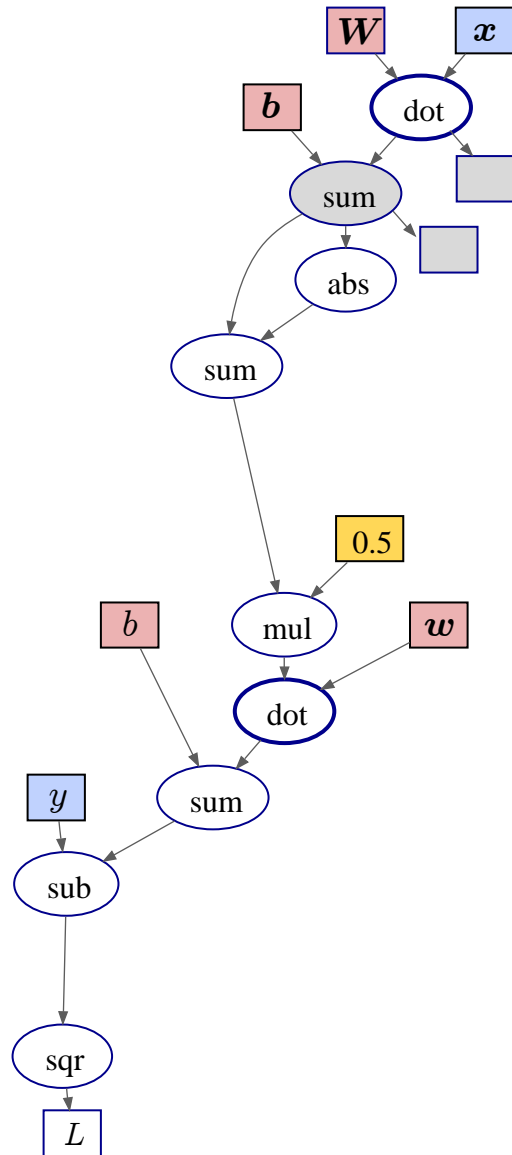
## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Item-wise loss function, FF neural network with ReLU as non-linearity



# An aside: Flow Graph



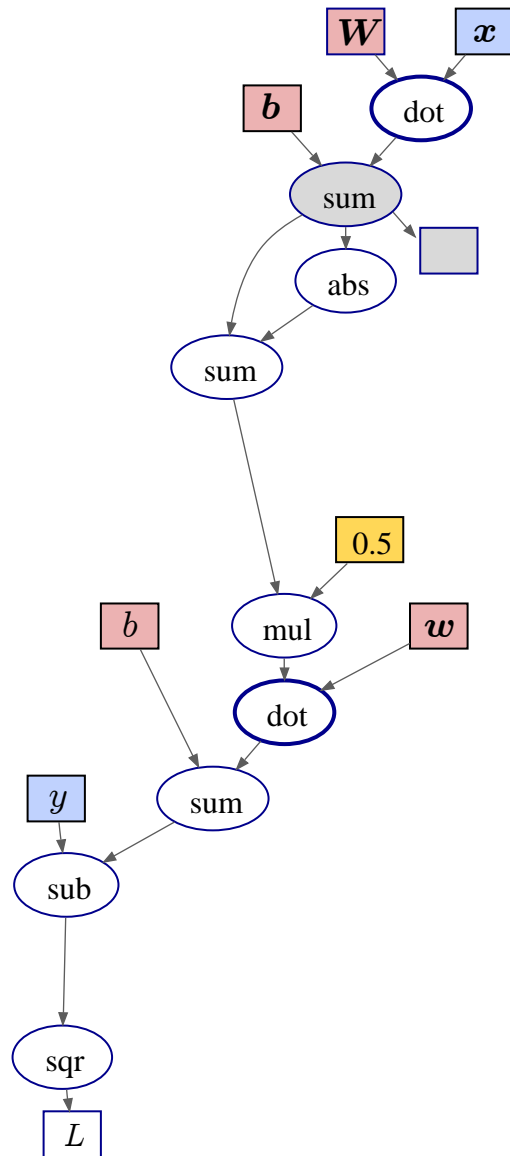
## ■ Computing the Flow Graph

*This is no longer necessary*

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

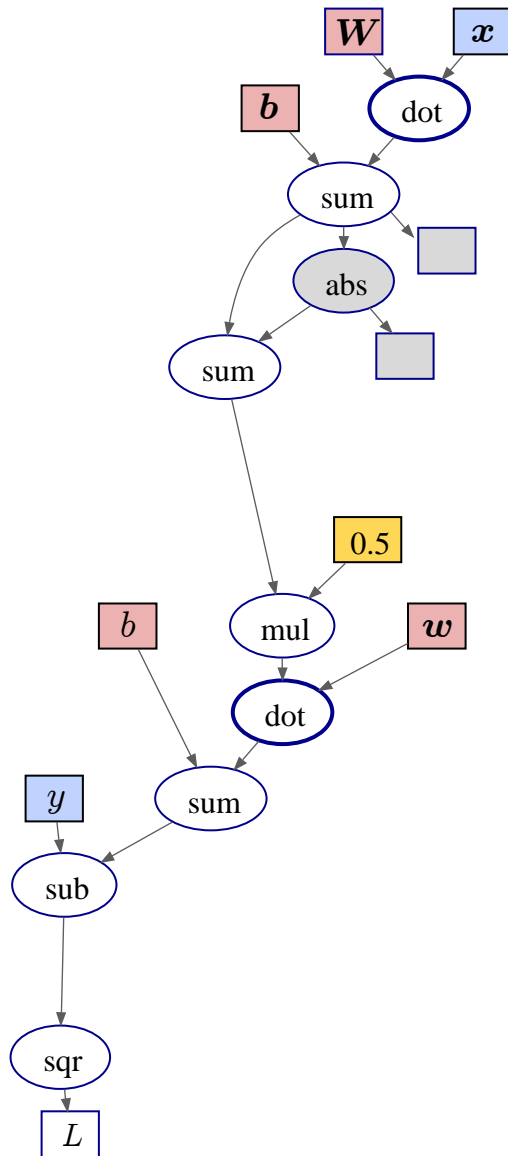


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

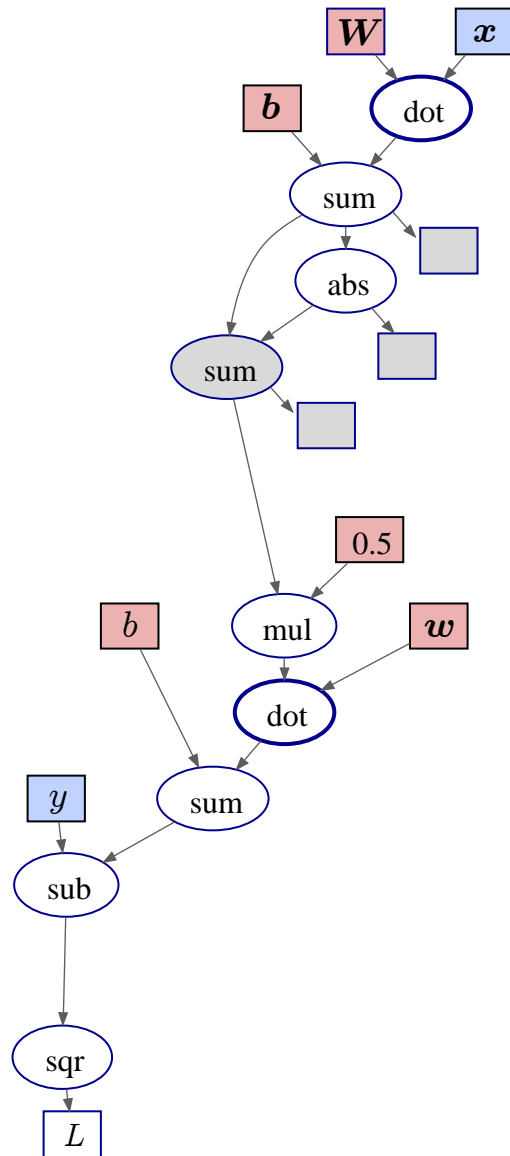


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

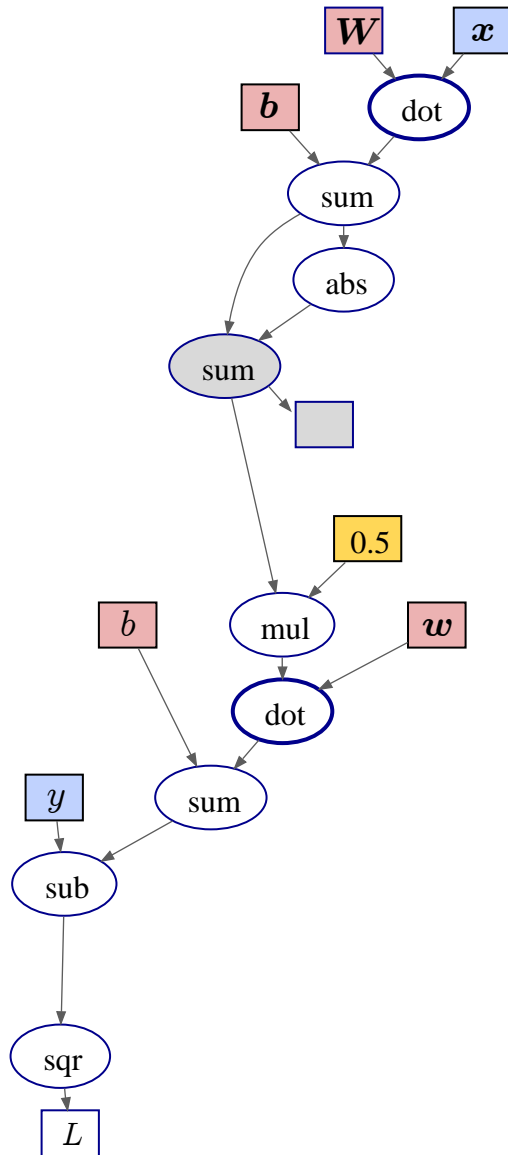


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

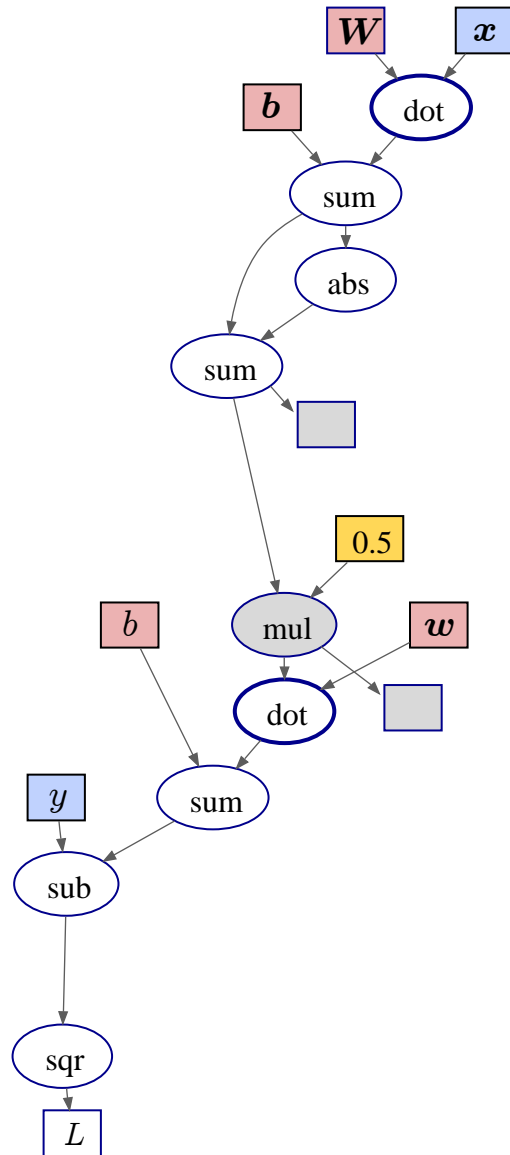


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

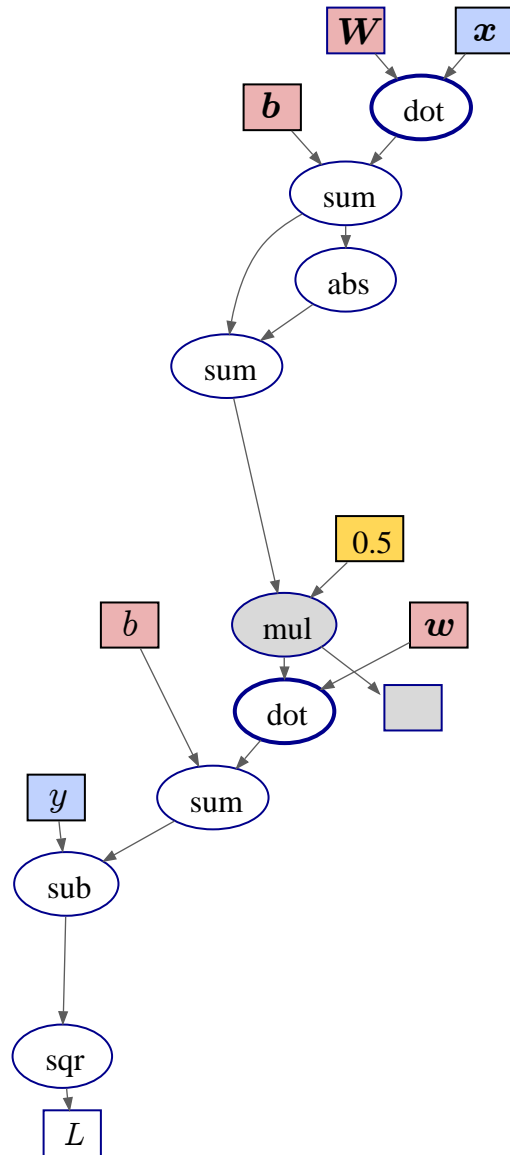


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

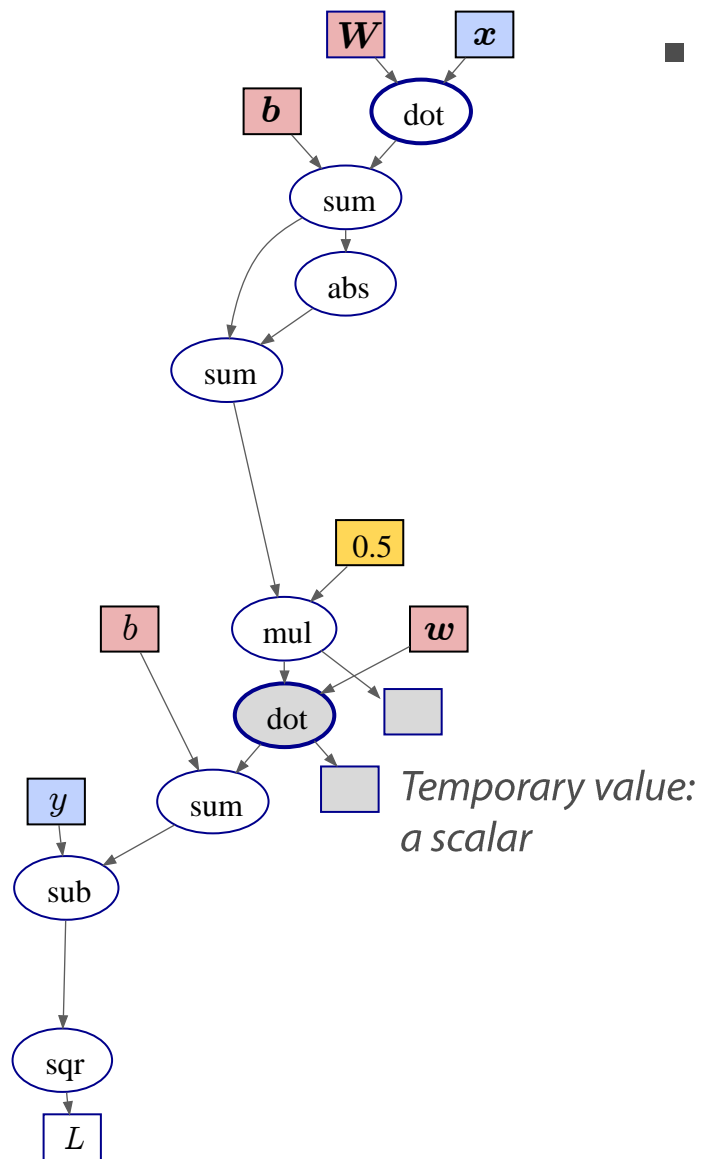


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph



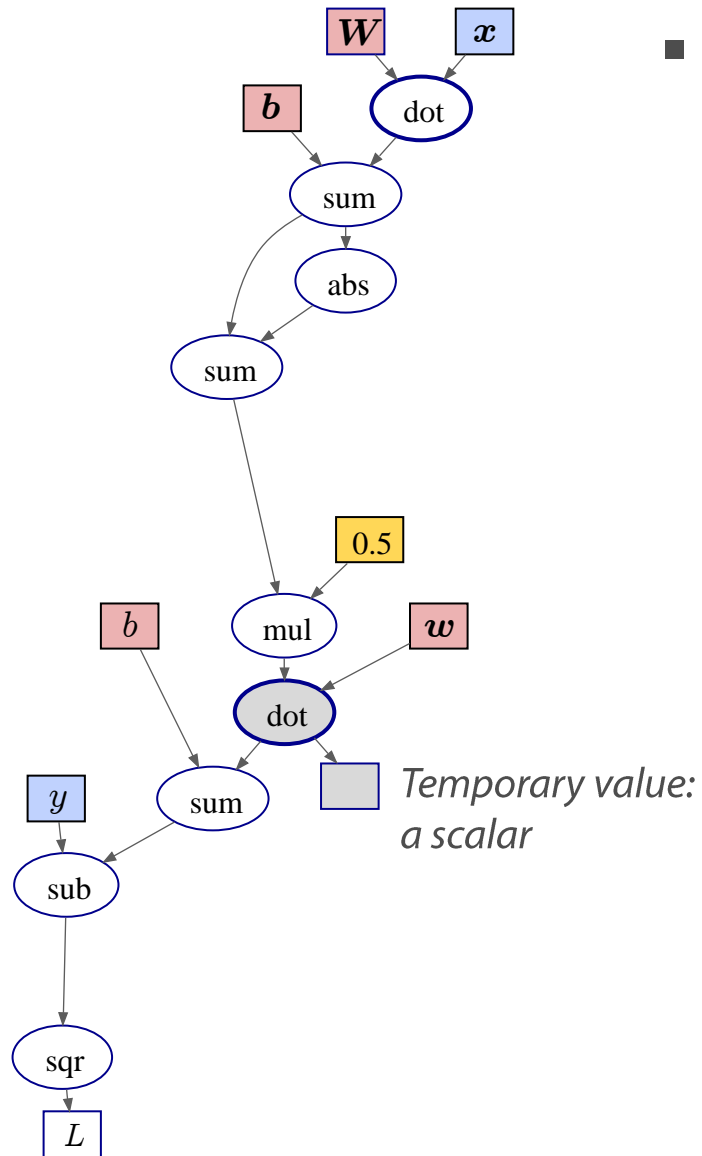
## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*



# An aside: Flow Graph

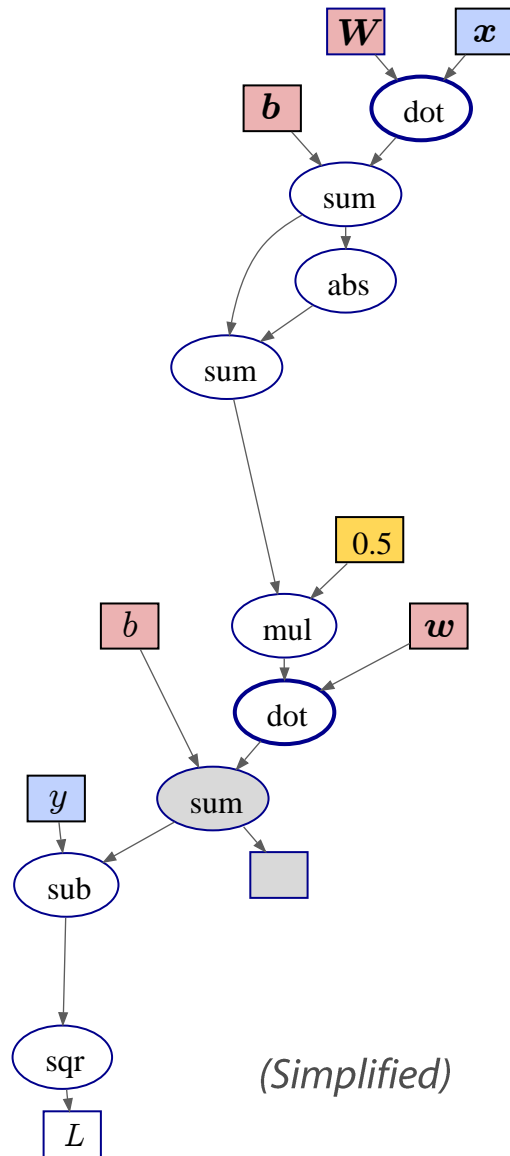


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph

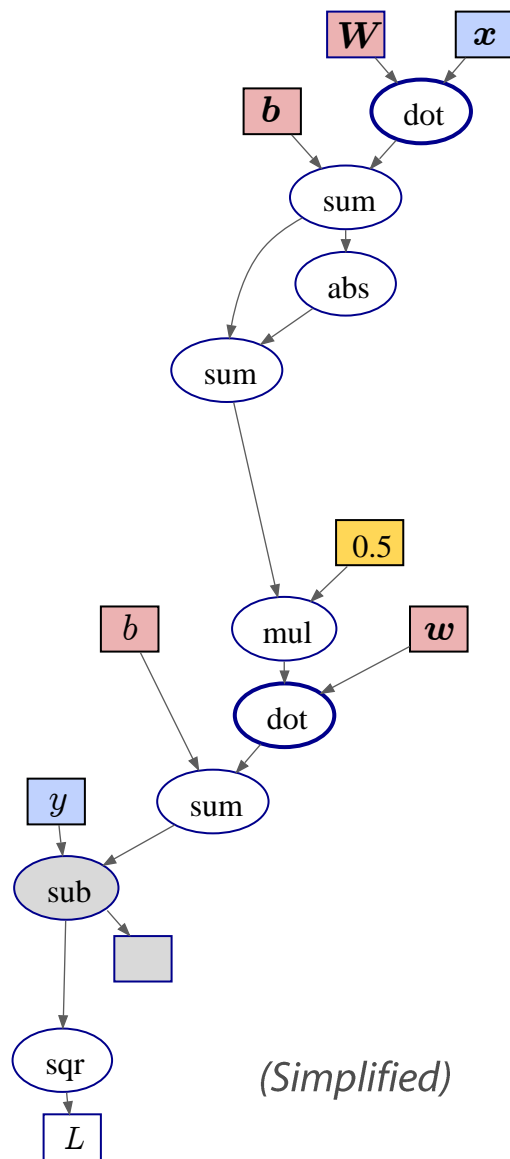


## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph



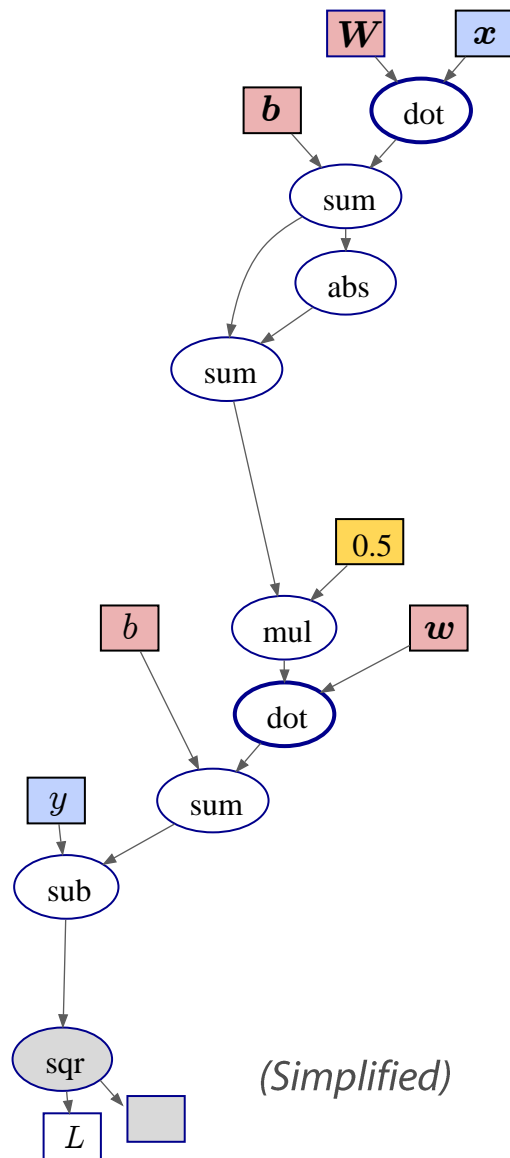
(Simplified)

## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

# An aside: Flow Graph



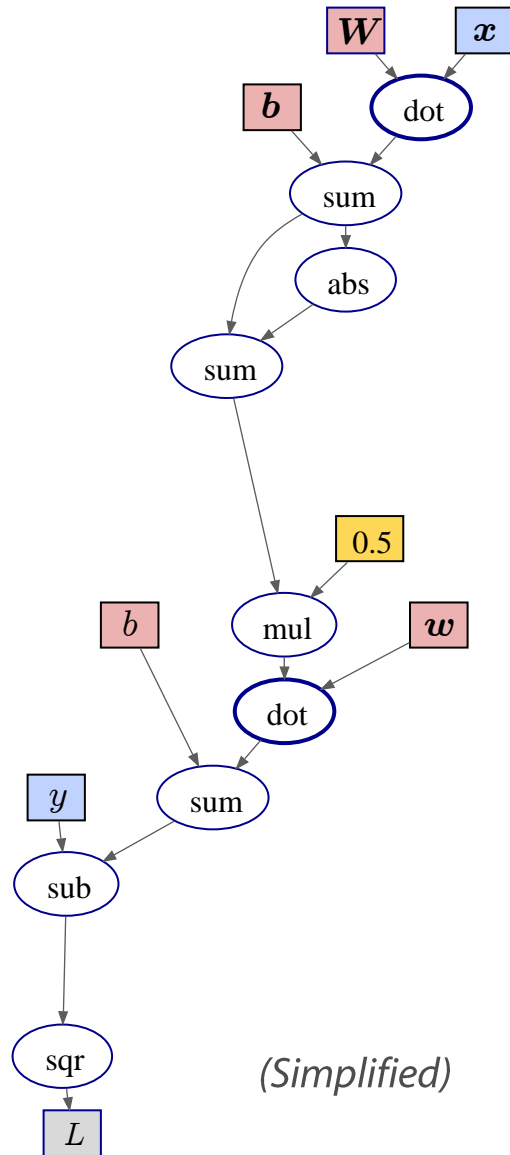
(Simplified)

## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Item-wise loss function, FF neural network with ReLU as non-linearity

# An aside: Flow Graph



(Simplified)

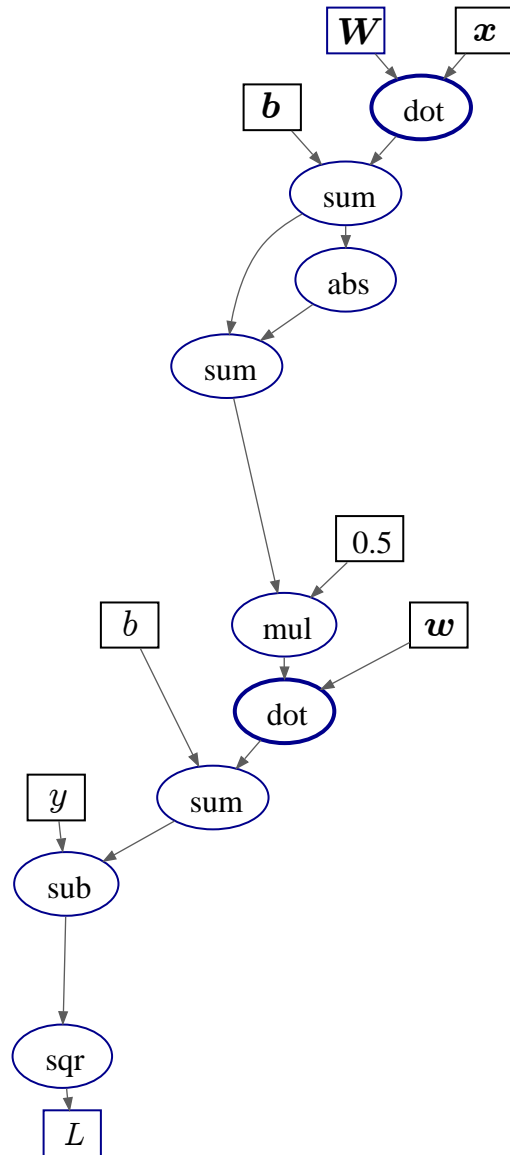
## ■ Computing the Flow Graph

$$L(\tilde{y}, y) = (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*Item-wise loss function, FF neural network with ReLU as non-linearity*

*Autodiff:  
Automatic Differentiation  
of Flow Graphs*

# Computing Gradients

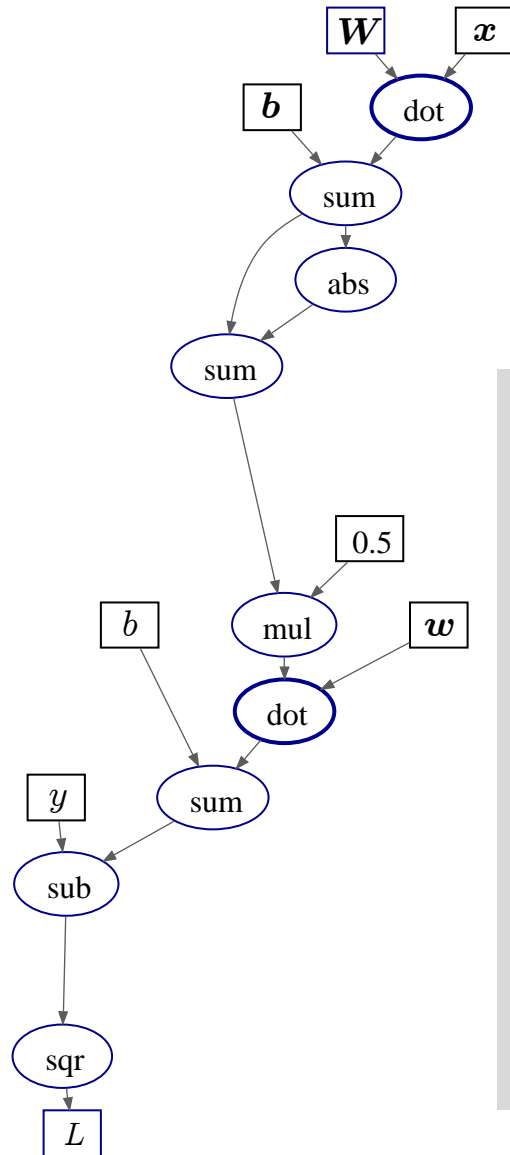


- **Computing one gradient of the flow graph**

$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*This is the gradient we want to compute  
(remember this is just one of the four)*

# Computing Gradients



- **Computing one gradient of the flow graph**

$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

*This is the gradient we want to compute  
(remember this is just one of the four)*

Chain rule for derivatives (*single argument*)

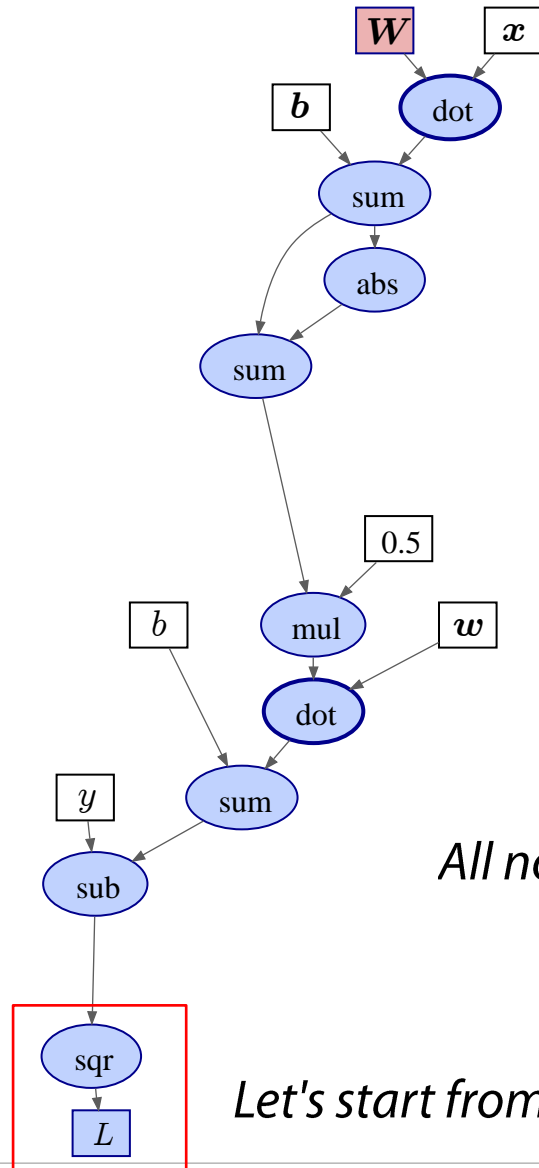
$$\frac{\partial}{\partial \boldsymbol{\vartheta}} f(g(\boldsymbol{\vartheta})) = \frac{\partial}{\partial g(\boldsymbol{\vartheta})} f(g(\boldsymbol{\vartheta})) \frac{\partial}{\partial \boldsymbol{\vartheta}} g(\boldsymbol{\vartheta})$$

Chain rule for derivatives (*multiple arguments*)

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} f(g(\boldsymbol{\vartheta}), h(\boldsymbol{\vartheta})) = \frac{\partial}{\partial g(\boldsymbol{\vartheta})} f(g(\boldsymbol{\vartheta}), h(\boldsymbol{\vartheta})) \frac{\partial}{\partial \boldsymbol{\vartheta}} g(\boldsymbol{\vartheta}) + \frac{\partial}{\partial h(\boldsymbol{\vartheta})} f(g(\boldsymbol{\vartheta}), h(\boldsymbol{\vartheta})) \frac{\partial}{\partial \boldsymbol{\vartheta}} h(\boldsymbol{\vartheta})$$



# Computing Gradients

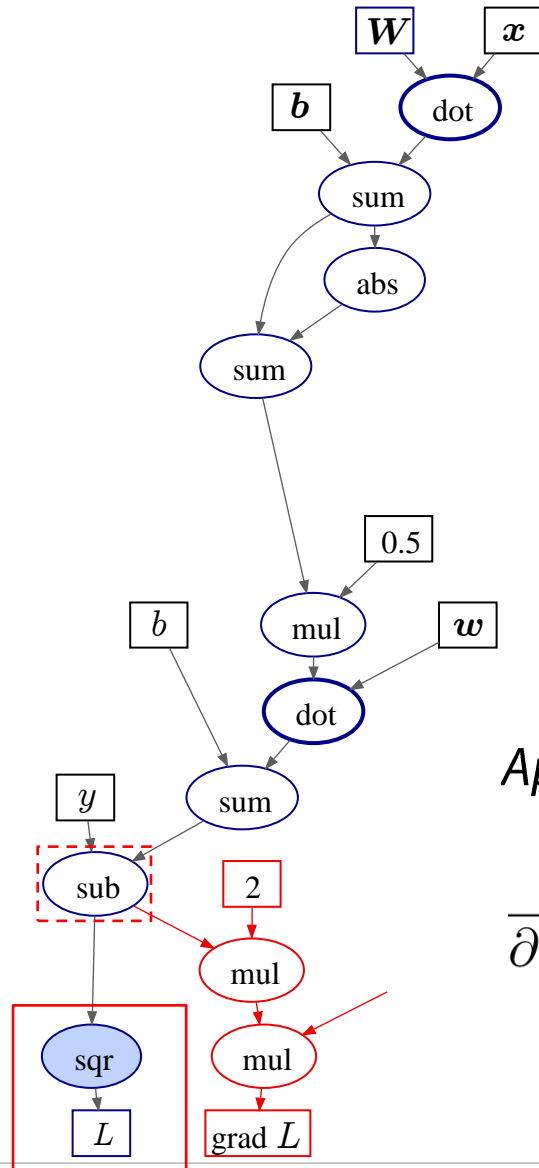


$$\frac{\partial}{\partial \mathbf{W}} (w \cdot \text{ReLU}(\mathbf{W}x + \mathbf{b}) + b - y)^2$$

All nodes depending on  $\mathbf{W}$  are marked in blue

Let's start from here (i.e. **backpropagation**, a.k.a. **reverse accumulation**)

# Computing Gradients

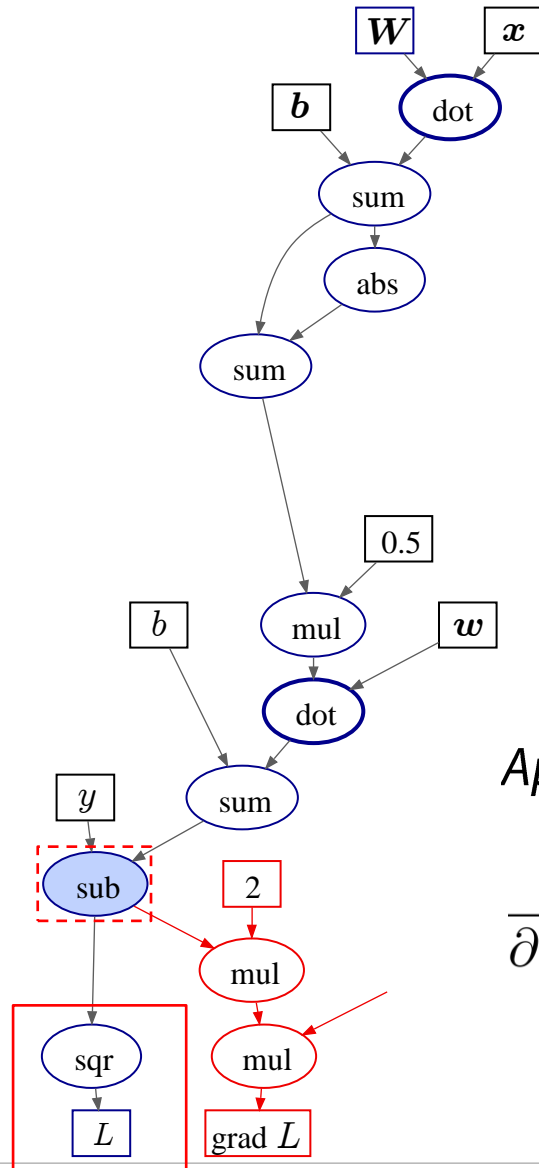


$$\frac{\partial}{\partial \mathbf{W}} (w \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{b} - \mathbf{y})^2$$

Apply the chain rule to the sqr node

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W})^2 &= \frac{\partial}{\partial f(\mathbf{W})} f(\mathbf{W})^2 \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= 2 \cdot f(\mathbf{W}) \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$

# Computing Gradients

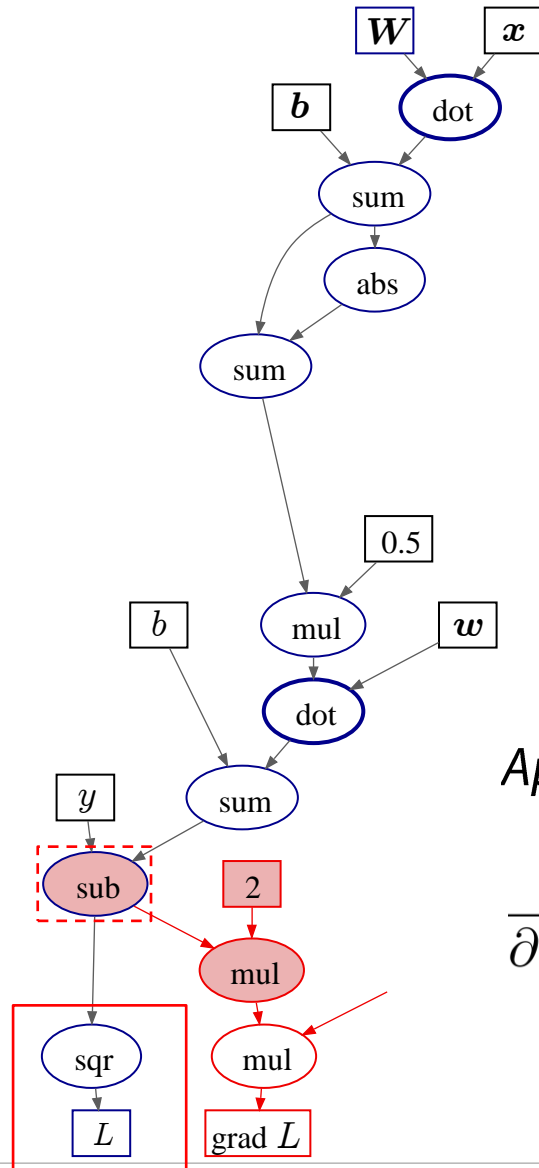


$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Apply the chain rule to the sqr node

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W})^2 &= \frac{\partial}{\partial f(\mathbf{W})} f(\mathbf{W})^2 \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= 2 \cdot f(\mathbf{W}) \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$

# Computing Gradients

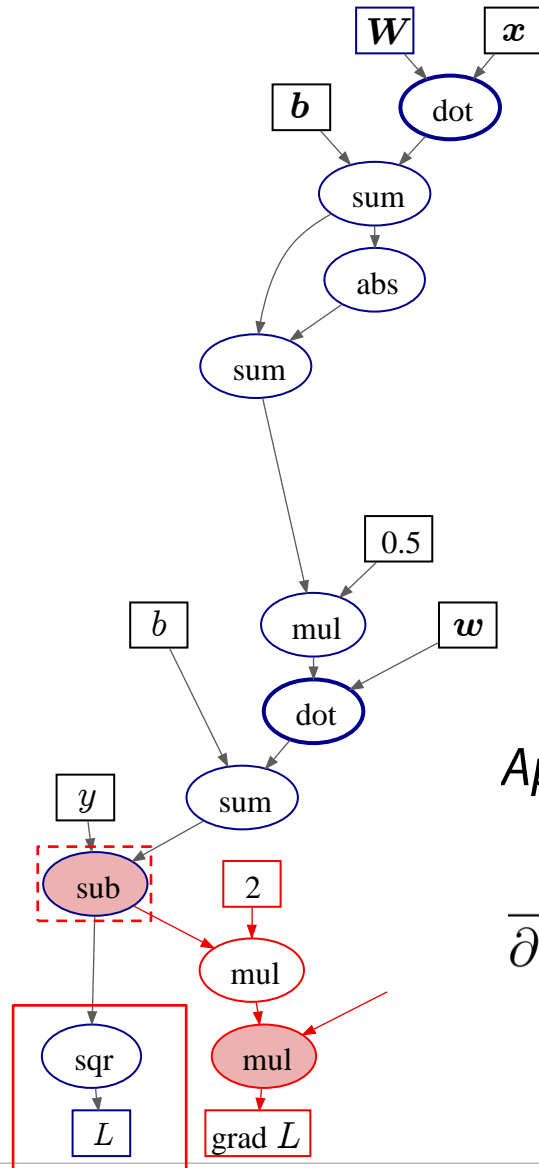


$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Apply the chain rule to the sqr node

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W})^2 &= \frac{\partial}{\partial f(\mathbf{W})} f(\mathbf{W})^2 \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= 2 \cdot f(\mathbf{W}) \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$

# Computing Gradients



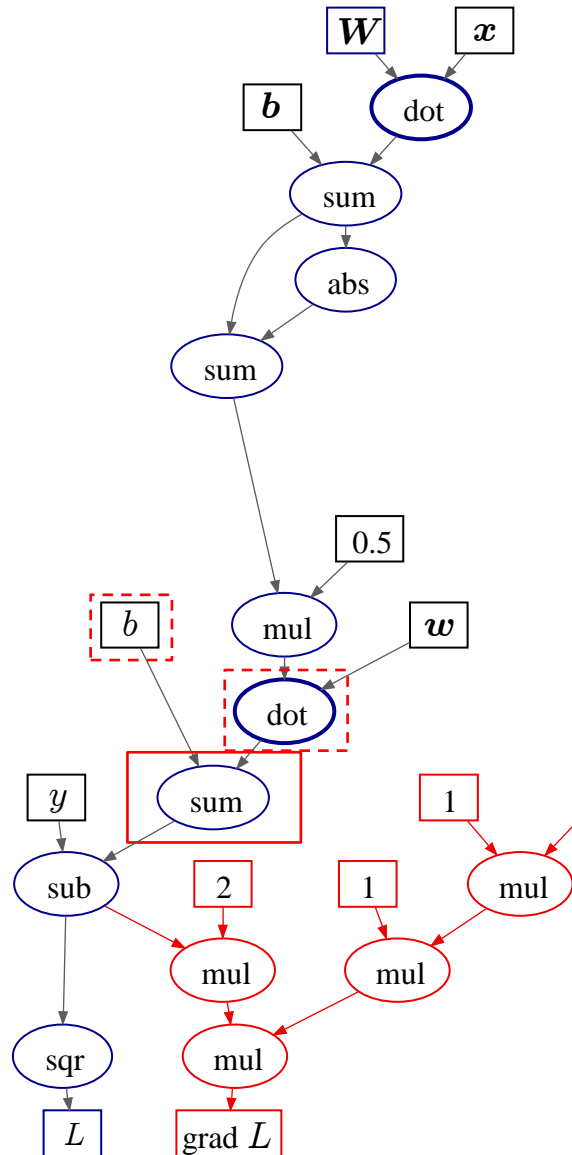
$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Apply the chain rule to the sqr node

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W})^2 &= \frac{\partial}{\partial f(\mathbf{W})} f(\mathbf{W})^2 \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= 2 \cdot f(\mathbf{W}) \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$



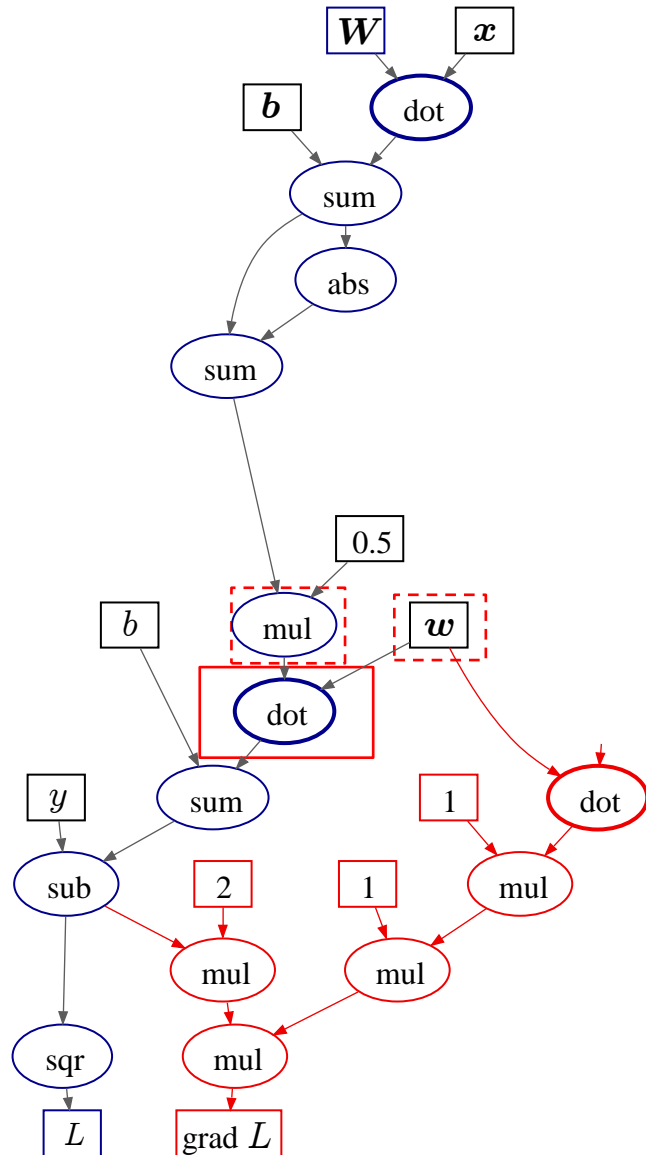
# Computing Gradients



$$\frac{\partial}{\partial \mathbf{W}} (w \cdot \text{ReLU}(\mathbf{W}x + \mathbf{b}) + b - y)^2$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} (f(\mathbf{W}) + b) &= \frac{\partial}{\partial f(\mathbf{W})} (f(\mathbf{W}) + b) \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= 1 \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$

# Computing Gradients

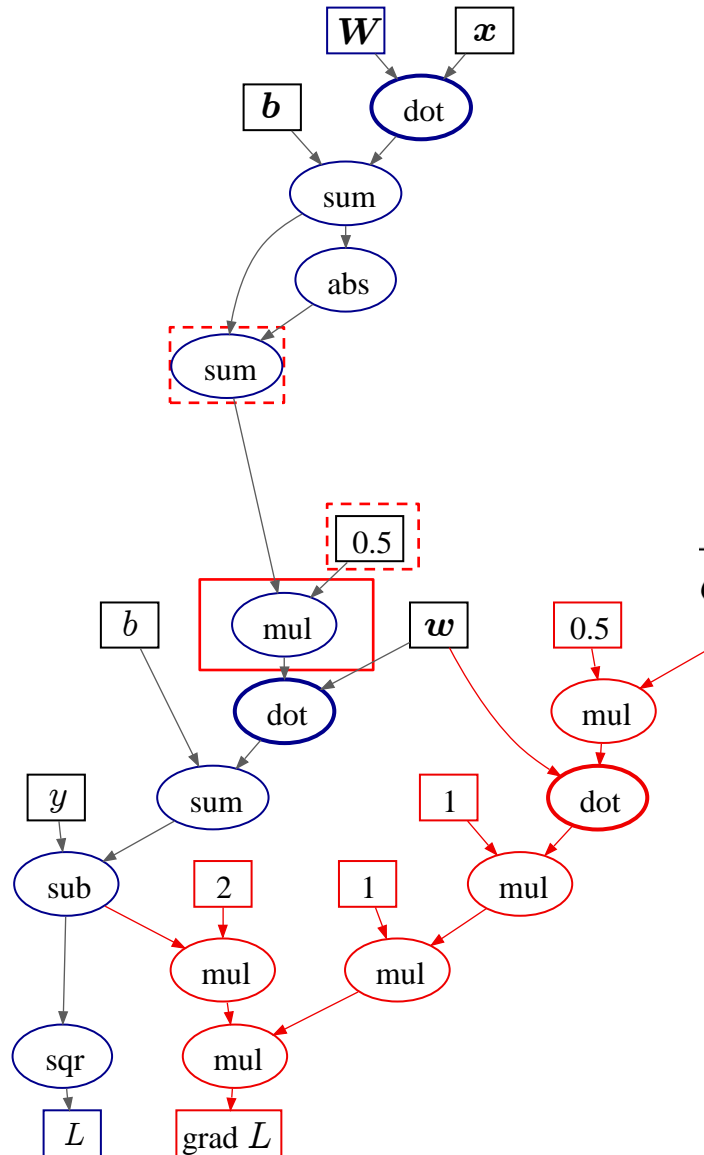


$$\frac{\partial}{\partial \mathbf{W}} (w \cdot \text{ReLU}(\mathbf{W}x + \mathbf{b}) + b - y)^2$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} (w \cdot f(\mathbf{W})) &= \frac{\partial}{\partial f(\mathbf{W})} (w \cdot f(\mathbf{W})) \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= w \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$



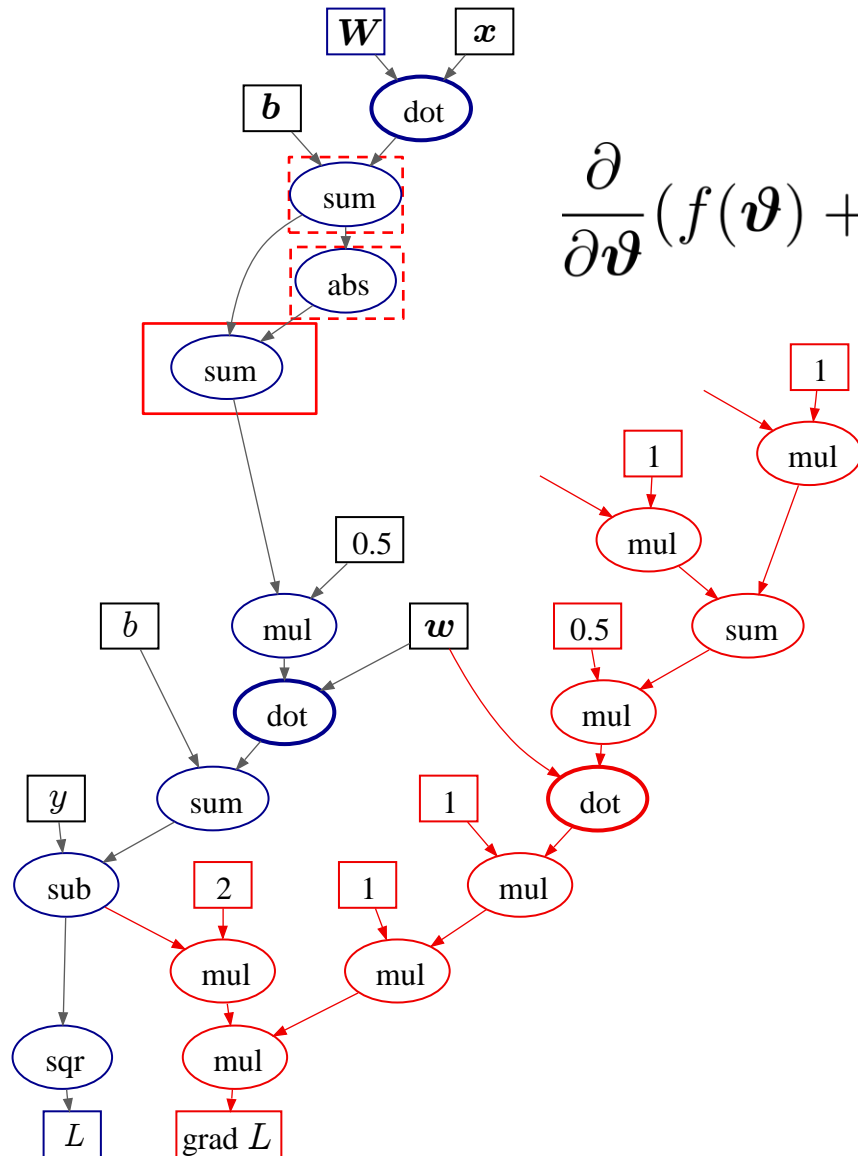
# Computing Gradients



$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

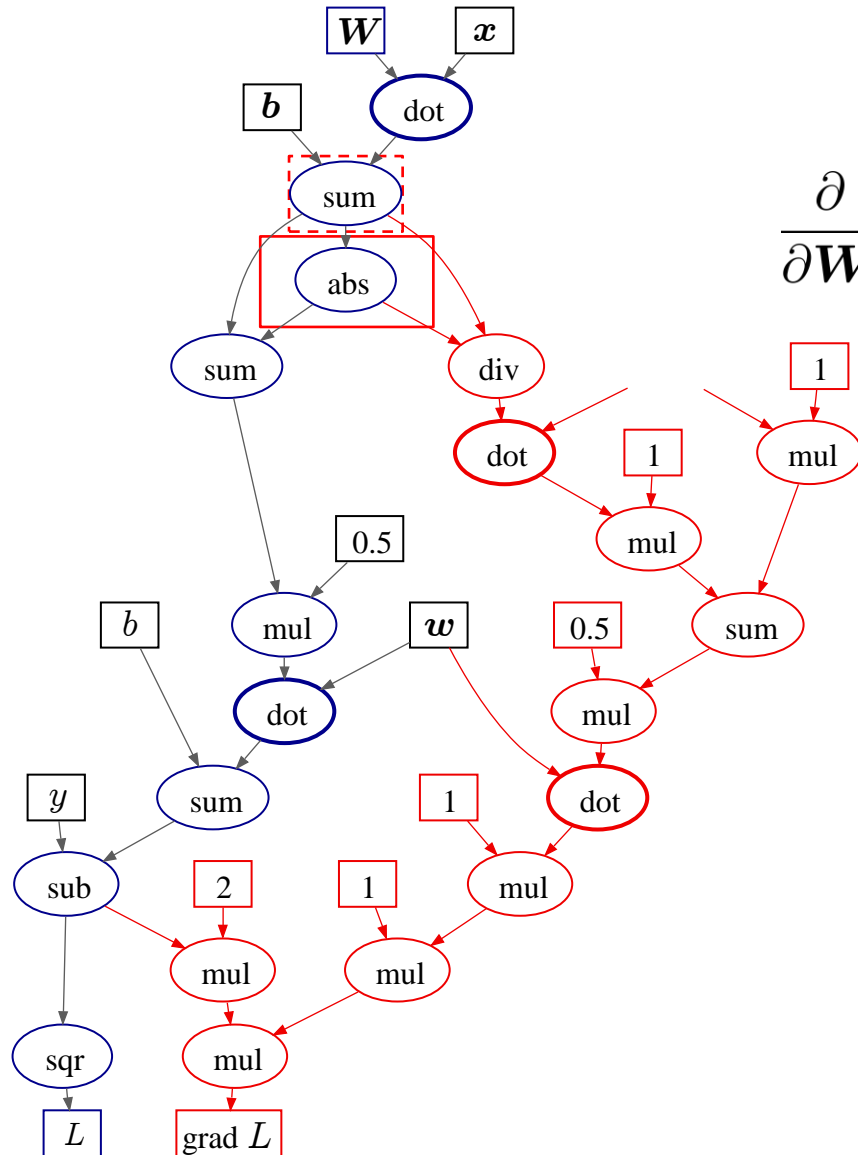
$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} (0.5 \cdot f(\mathbf{W})) &= \frac{\partial}{\partial f(\mathbf{W})} (0.5 \cdot f(\mathbf{W})) \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= 0.5 \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$

# Computing Gradients



$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} (f(\boldsymbol{\vartheta}) + g(\boldsymbol{\vartheta})) &= \frac{\partial}{\partial f(\boldsymbol{\vartheta})} (f(\boldsymbol{\vartheta}) + g(\boldsymbol{\vartheta})) \frac{\partial}{\partial \boldsymbol{\vartheta}} f(\boldsymbol{\vartheta}) \\ &+ \frac{\partial}{\partial g(\boldsymbol{\vartheta})} (f(\boldsymbol{\vartheta}) + g(\boldsymbol{\vartheta})) \frac{\partial}{\partial \boldsymbol{\vartheta}} g(\boldsymbol{\vartheta}) \\ &= 1 \cdot \frac{\partial}{\partial \boldsymbol{\vartheta}} f(\boldsymbol{\vartheta}) + 1 \cdot \frac{\partial}{\partial \boldsymbol{\vartheta}} g(\boldsymbol{\vartheta}) \end{aligned}$$

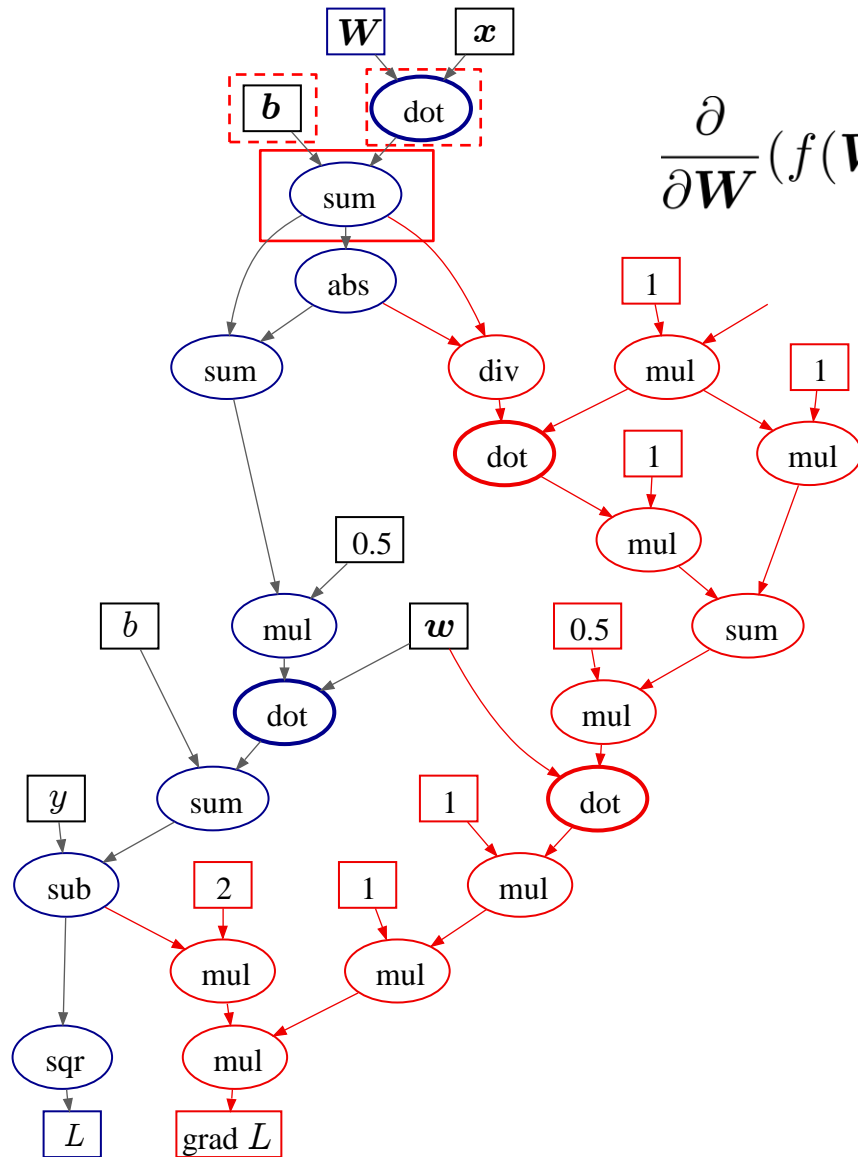
# Computing Gradients



$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} |f(\mathbf{W})| &= \frac{\partial}{\partial f(\mathbf{W})} |f(\mathbf{W})| \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \\ &= \frac{f(\mathbf{W})}{|f(\mathbf{W})|} \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W}) \end{aligned}$$

Clearly, this term is not defined for any  $W_{ij} = 0$   
 (Typically, this is a protected division  $\frac{x}{0} := 1$ )

# Computing Gradients

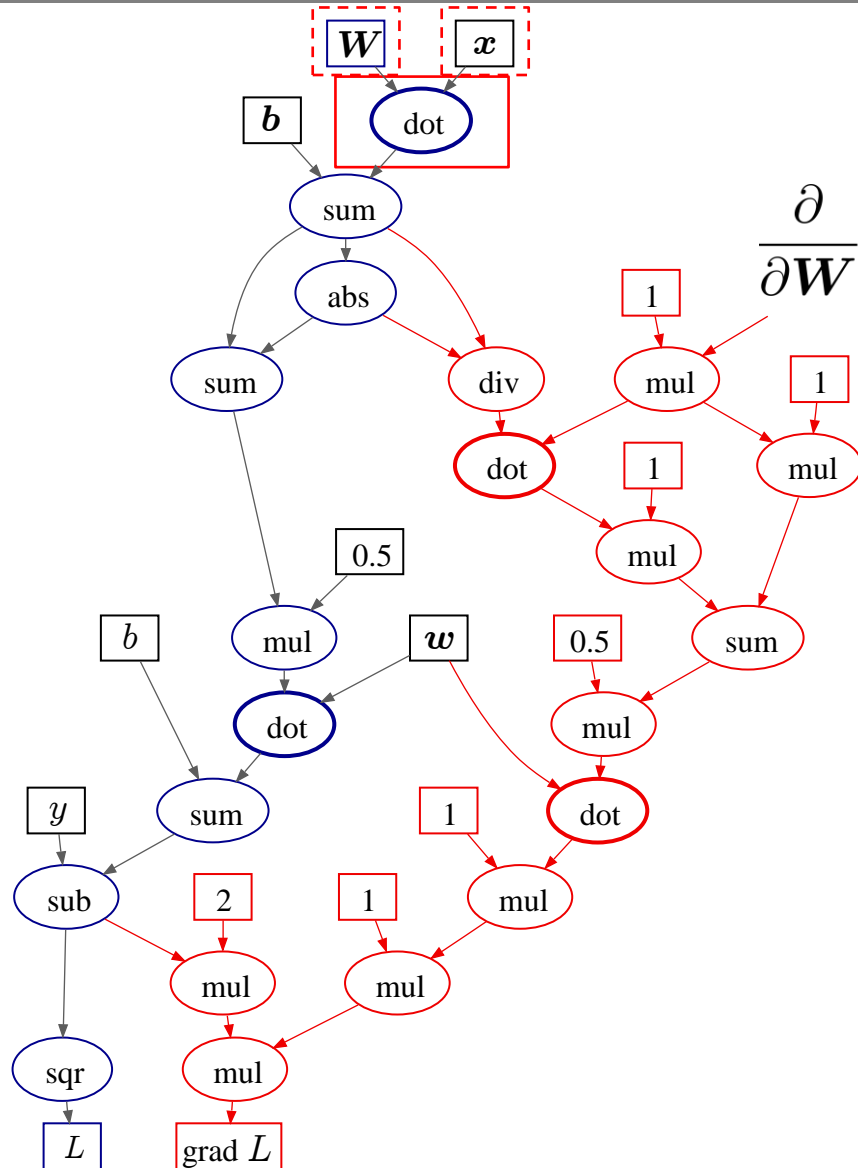


$$\frac{\partial}{\partial \mathbf{W}} (f(\mathbf{W}) + b) = \frac{\partial}{\partial f(\mathbf{W})} (f(\mathbf{W}) + b) \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W})$$

$$= 1 \cdot \frac{\partial}{\partial \mathbf{W}} f(\mathbf{W})$$



# Computing Gradients



$$\frac{\partial}{\partial W} (W \cdot x)$$

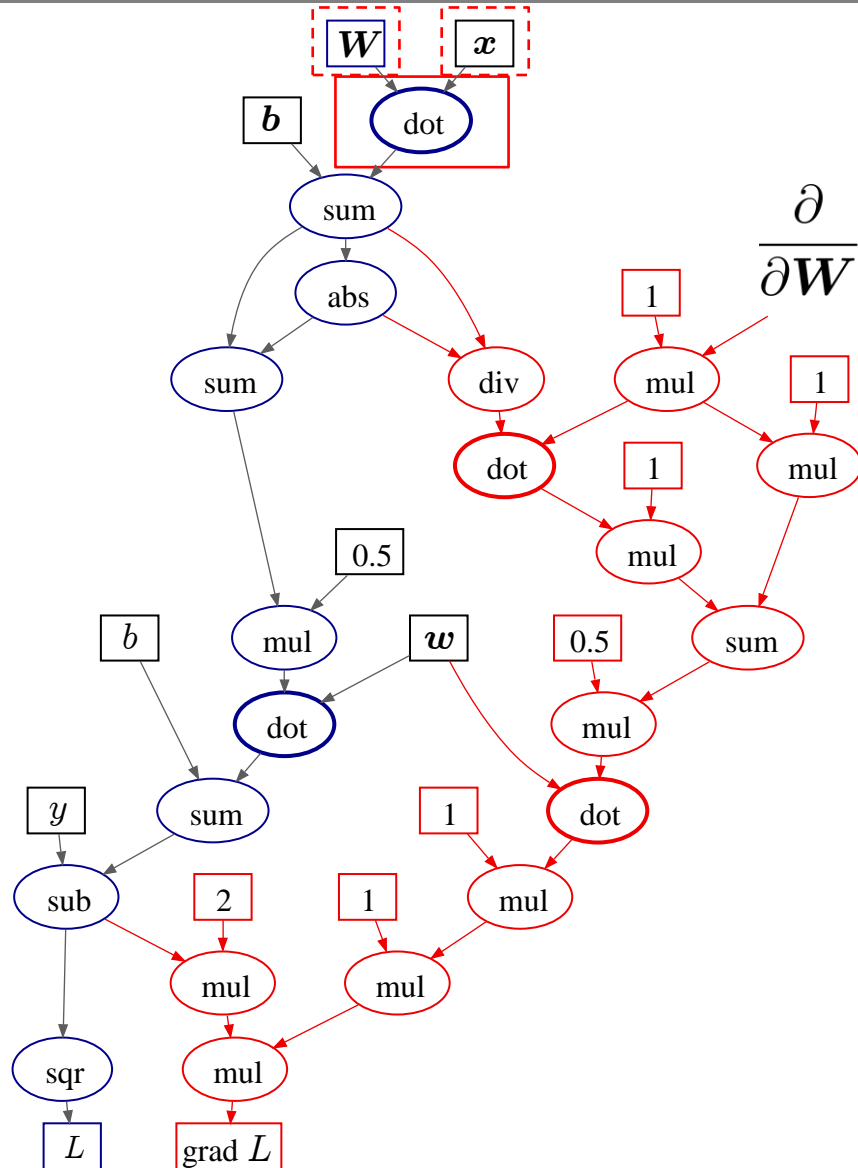
This is a third-order tensor

$$\left( \frac{\partial}{\partial W} (W \cdot x) \right)_{ijk} = \frac{\partial}{\partial W_{kj}} (W \cdot x)_i$$

Its  $ijk$ -th component

Note the inversion of indices

# Computing Gradients



$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{W} \cdot \mathbf{x})$$

This is a third-order tensor

$$\left( \frac{\partial}{\partial \mathbf{W}} (\mathbf{W} \cdot \mathbf{x}) \right)_{ijk} = \frac{\partial}{\partial W_{kj}} (\mathbf{W} \cdot \mathbf{x})_i$$

Its  $ijk$ -th component

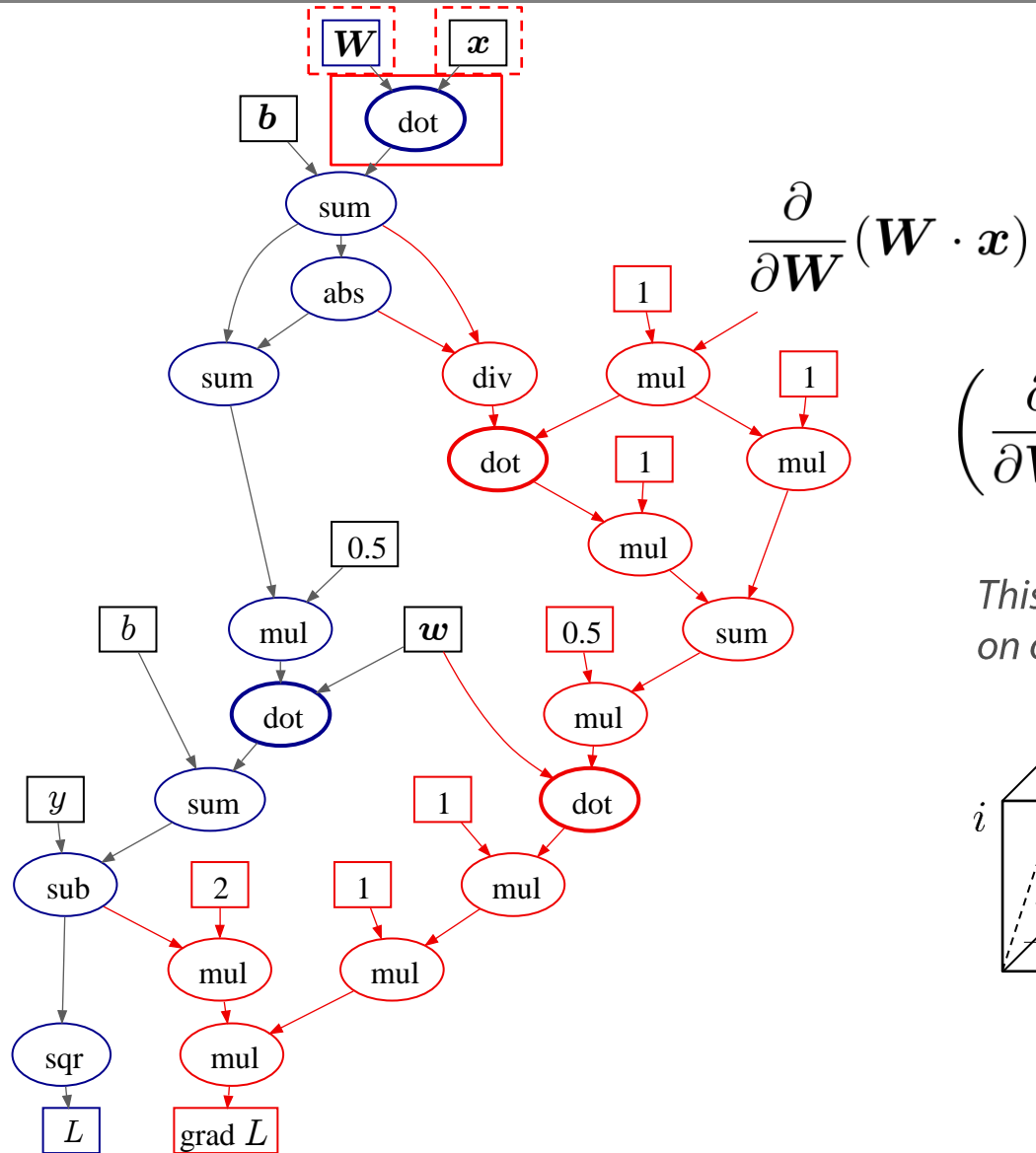
Note the inversion of indices

$$= \frac{\partial}{\partial W_{kj}} (\mathbf{W}_{i,:} \cdot \mathbf{x})$$

The  $i$ -th line in the matrix

$$= \begin{cases} 0 & k \neq i \\ x_j & k = i \end{cases}$$

# Computing Gradients

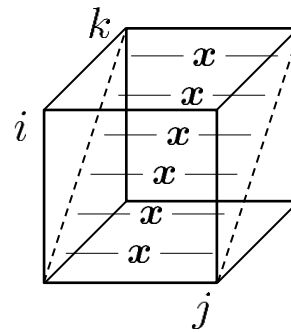


$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{W} \cdot \mathbf{x})$$

Putting it all together...

$$\left( \frac{\partial}{\partial \mathbf{W}} (\mathbf{W} \cdot \mathbf{x}) \right)_{ijk} = \begin{cases} 0 & k \neq i \\ x_j & k = i \end{cases}$$

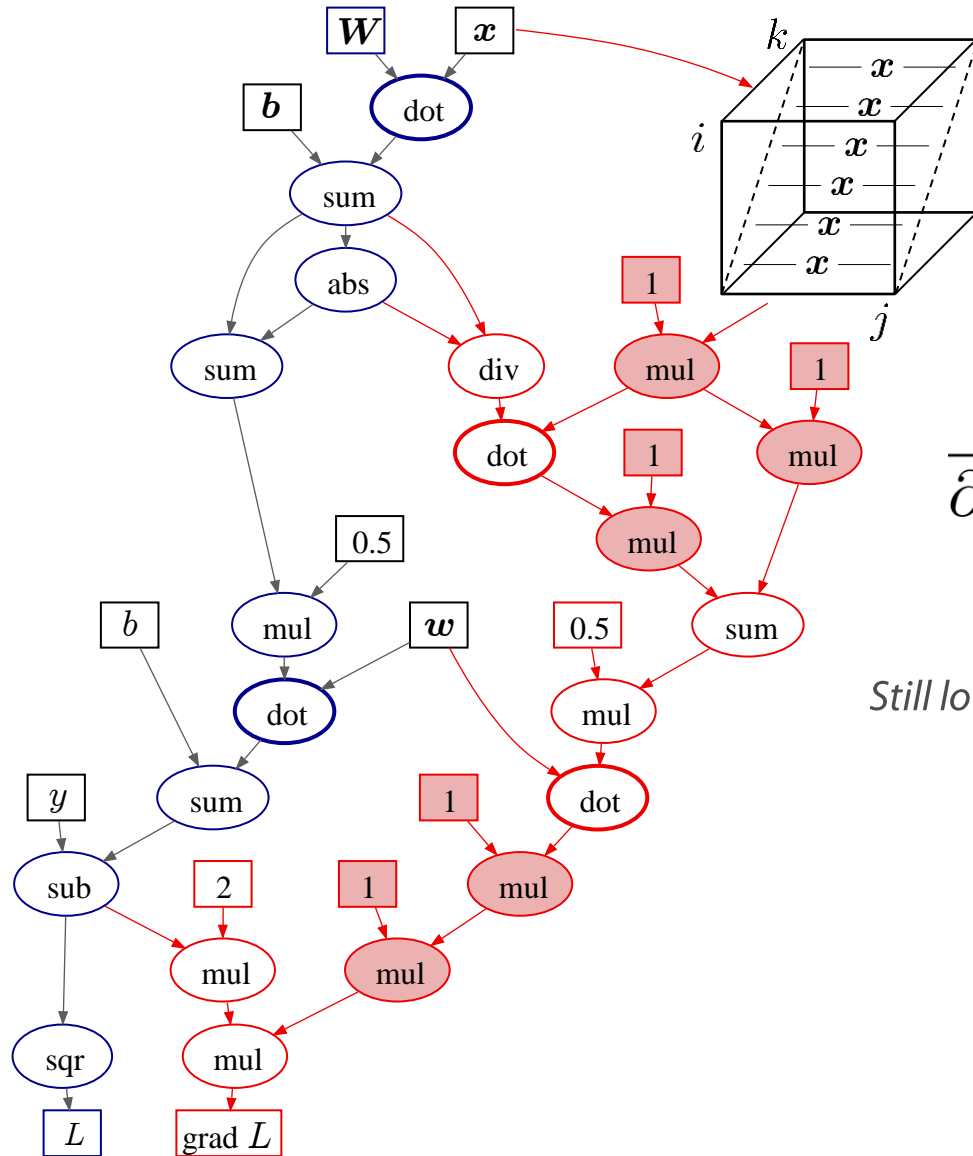
This 'thing' is a cube having copies of  $\mathbf{x}$  on one diagonal 'plane' and zeros elsewhere







# Computing Gradients

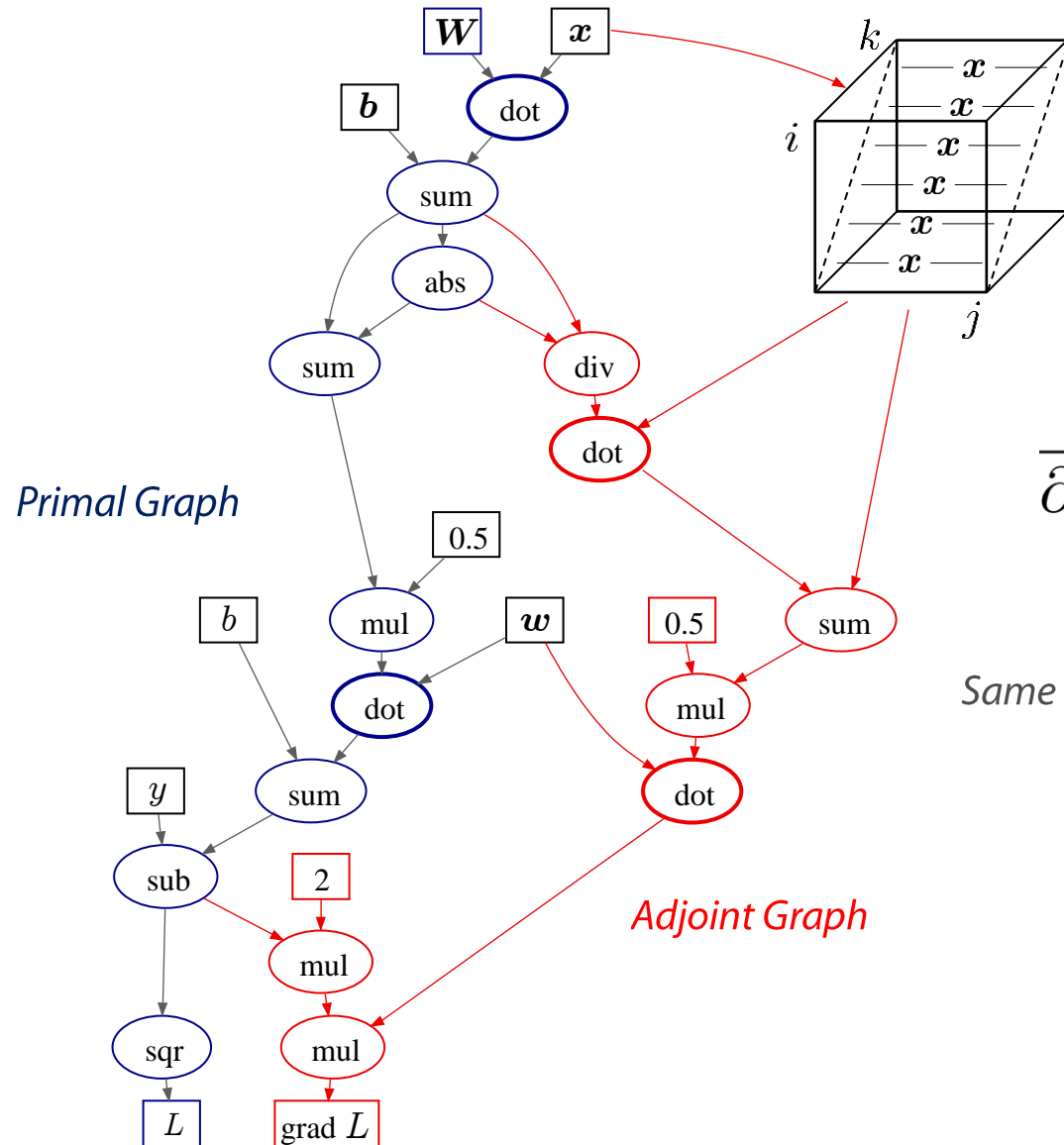


The representation of this can be optimized too

$$\frac{\partial}{\partial \mathbf{W}} (\mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b - y)^2$$

Still lots of useless operations

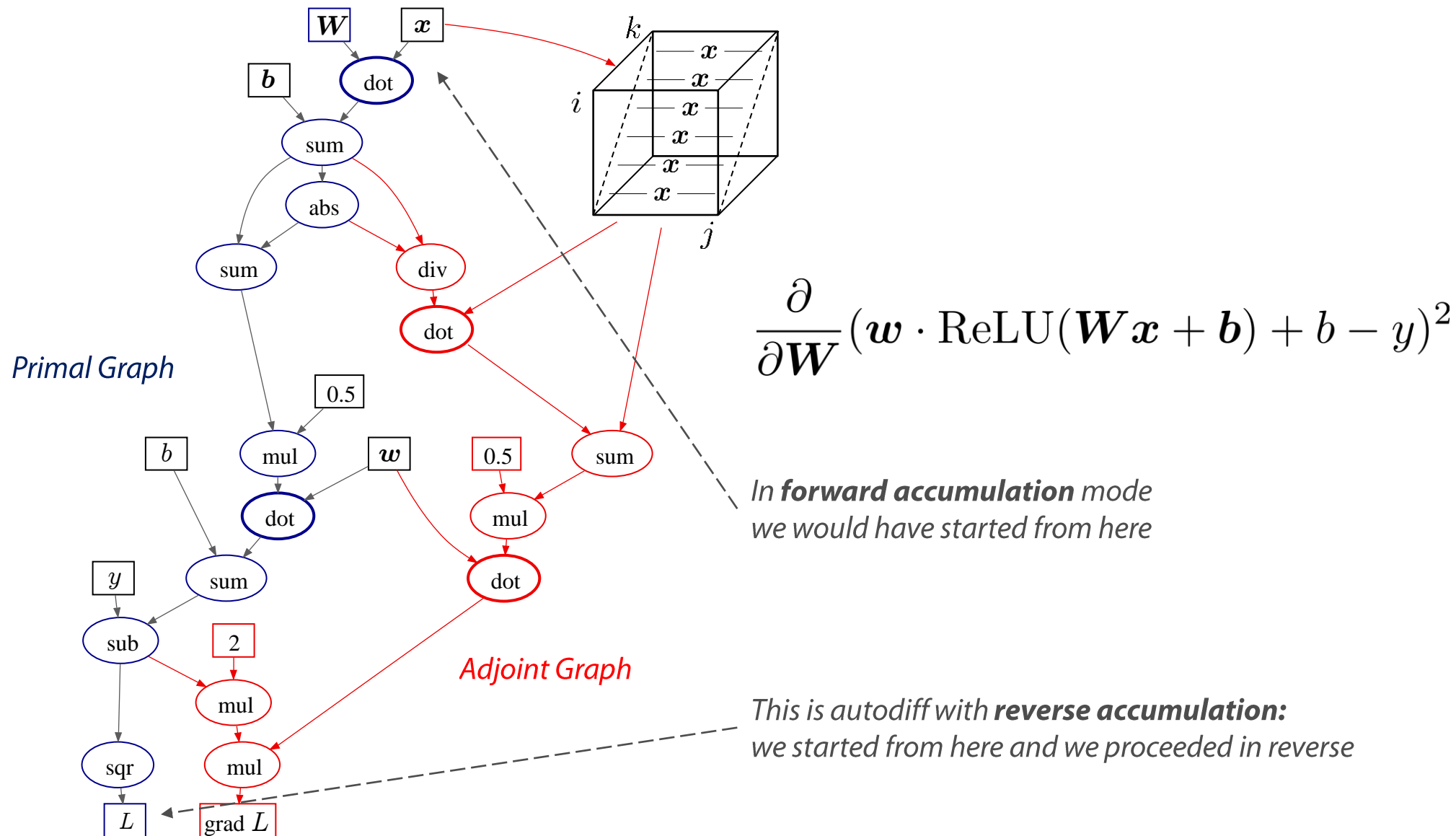
# Computing Gradients



$$\frac{\partial}{\partial W} (w \cdot \text{ReLU}(Wx + b) + b - y)^2$$

Same graph, after some pruning

# Computing Gradients



# *(Mini) Batches in Matrix Form*

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Let's focus first on  $\mathbf{W}\mathbf{x}$

by defining  $\mathbf{X} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix}$  input data in matrix form (**item index first**)

Then we can write

$$\mathbf{W}\mathbf{X}^T = \begin{bmatrix} \mathbf{W}\mathbf{x}^{(1)} & \dots & \mathbf{W}\mathbf{x}^{(N)} \end{bmatrix}$$

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Consider then  $(\mathbf{W}\mathbf{x} + \mathbf{b})$

by defining

$$\hat{\mathbf{X}} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} & 1 \end{bmatrix} \quad \hat{\mathbf{W}} := \begin{bmatrix} \mathbf{W} & | \\ \mathbf{b} & | \end{bmatrix}$$

Then we could write

$$\hat{\mathbf{W}}\hat{\mathbf{X}}^T = \begin{bmatrix} \mathbf{W}\mathbf{x}^{(1)} + \mathbf{b} & \dots & \mathbf{W}\mathbf{x}^{(N)} + \mathbf{b} \end{bmatrix}$$

Matrix  $\hat{\mathbf{W}}$  includes two parameters:  $\mathbf{W}$  and  $\mathbf{b}$   
*this may be inconvenient for autodiff..*

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Consider then  $(\mathbf{W}\mathbf{x} + \mathbf{b})$

and let's keep the definition

$$\mathbf{X} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix}$$

It could be convenient to redefine the operator  $+$  such that is interpreted as

$$\mathbf{W}\mathbf{X}^T + \mathbf{b} := \begin{bmatrix} \mathbf{W}\mathbf{x}^{(1)} & \dots & \mathbf{W}\mathbf{x}^{(N)} \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{b} & \dots & \mathbf{b} \end{bmatrix}}_{N \text{ times}}$$



# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Consider then  $(\mathbf{W}\mathbf{x} + \mathbf{b})$

and let's keep the definition

$$\mathbf{X} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix}$$

It could be convenient to redefine the operator  $+$  such that is interpreted as

$$\mathbf{W}\mathbf{X}^T + \mathbf{b} := \begin{bmatrix} \mathbf{W}\mathbf{x}^{(1)} & \dots & \mathbf{W}\mathbf{x}^{(N)} \\ | & & | \\ | & & | \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{b} & \dots & \mathbf{b} \\ | & & | \\ | & & | \end{bmatrix}}_{N \text{ times}}$$

*This is called **broadcasting***

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W} \mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W} \mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

where

$$\mathbf{X} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \quad \mathbf{y} := \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

—  
This is a matrix  $g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) \in \mathbb{R}^{h \times N}$   
(Note the **broadcast** with  $+$ )

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

This is a **row** vector

$$\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) = \mathbf{w}^T g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) \in \mathbb{R}^N$$

(The 'dot' operator **transposes vectors** automatically, as required)

**NOTE: automatic transposition applies to vectors only!**  
For any tensor beyond dimension 1, you need to do that on your own

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

—  
This is also a **row** vector  $\in \mathbb{R}^N$ , after a **broadcast** on  $b$

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

This is also a **row** vector  $\in \mathbb{R}^N$ , after a **broadcast** on  $b$

... whereas this is a **column** vector  $\in \mathbb{R}^N$

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

This is also a **row** vector  $\in \mathbb{R}^N$ , after a **broadcast** on  $b$

... whereas this is a **column** vector  $\in \mathbb{R}^N$

(Also, the  $-$  operator **transposes** vectors automatically, as required)

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

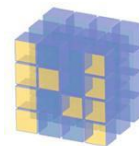
$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

This is also a **row** vector  $\in \mathbb{R}^N$ , after a **broadcast** on  $b$

... whereas this is a **column** vector  $\in \mathbb{R}^N$

(Also, the  $-$  operator **transposes** vectors automatically, as required)

A similar behavior of operators is standard in



NumPy



TensorFlow





# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{X}^T + \mathbf{b}) + b) - \mathbf{y})^2$$

This is a matrix  $\mathbf{W}\mathbf{X}^T \in \mathbb{R}^{h \times N}$

*Ouch! No item index first ...*

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((g(\mathbf{X}\mathbf{W}^T + \mathbf{b}) \cdot \mathbf{w} + b) - \mathbf{y})^2$$

This is a matrix  $\mathbf{X}\mathbf{W}^T \in \mathbb{R}^{N \times h}$

*Item index first!*

# More on Matrix Forms

Say it with matrices...

We may want to get rid of the summation when computing the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

Using broadcasting operators, we can express the above as

$$L(D) = \frac{1}{N} ((g(\mathbf{X}\mathbf{W}^T + \mathbf{b}) \cdot \mathbf{w} + b) - \mathbf{y})^2$$

This is a matrix  $\mathbf{X}\mathbf{W}^T \in \mathbb{R}^{N \times h}$

This is a **column** vector  $\in \mathbb{R}^h$   
(it will be transposed automatically)

**Item index first!**