

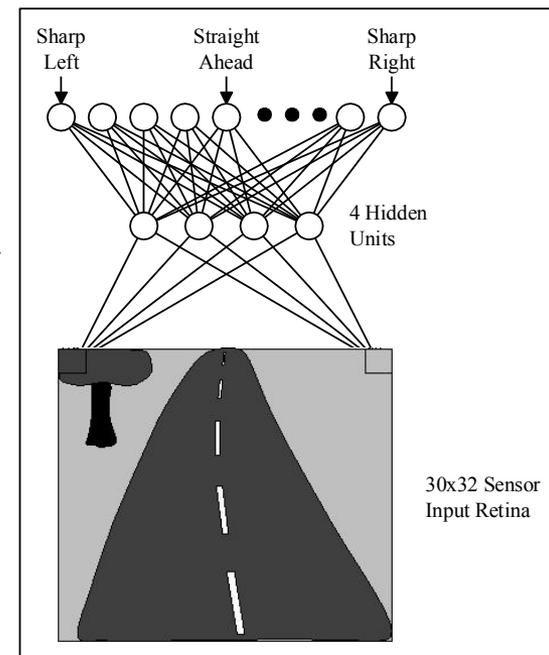


# **SGD and Deep Learning for Classification**

UCSB CS293S, 2020, T. Yang

# Motivation and Table of Content

- What we have learned so far for ranking and classification
  - Decision trees: entropy-based, or regression
  - Ensembles, boosting, and bagging. Random forests
- Focus of this slide set
  - Stochastic gradient descent (SGD) for general optimization
  - Derive weights for minimizing a loss function in a large network-based classification
  - Example of neural nets and optimization
- Why?
  - Successful in neural classification tasks for image and audio processing with machine learning
  - Effective for text oriented document classification and ranking



# Partial Derivatives and Gradient

## Single-variable functions

### Notation for the Derivative

$$\left. \begin{array}{l} f'(x) \\ y' \\ \frac{dy}{dx} \end{array} \right\} \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

## Multi-variable functions

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

## Gradient

Scalar-valued multivariable function

$$\nabla f(x_0, y_0, \dots) = \begin{bmatrix} \frac{\partial f}{\partial x}(x_0, y_0, \dots) \\ \frac{\partial f}{\partial y}(x_0, y_0, \dots) \\ \vdots \end{bmatrix}$$

Notation for gradient, called "nabla".

$\nabla f$  outputs a vector with all possible partial derivatives of  $f$ .

# Start with Simple Binary Text Classifier

Also called perceptron

$x$

$f(x)$

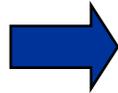
$y$

**Result classification:**

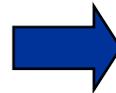
Positive, output +1

Negative, output -1

Hello,  
Do you want free printer  
cartridges? Why pay more  
when you can get them  
ABSOLUTELY FREE! Just

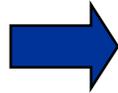


```
# free      : 2  
YOUR_NAME  : 0  
MISPELLED  : 2  
FROM_FRIEND : 0  
...
```

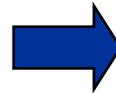


SPAM  
or  
+

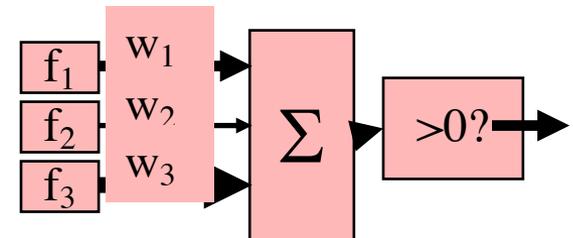
2



```
PIXEL-7,12 : 1  
PIXEL-7,13 : 0  
...  
NUM_LOOPS  : 1  
...
```



"2"



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

*Positive dot product  $w \cdot f$  means  
the positive class*

# SGD training for Binary Classifier

Figure out the weight vector from training instances

- Start with weights = 0
- For each training instance:
  - Classify with current weights
  - $f(x)$  is feature vector of  $x$

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

$$\begin{pmatrix} \# \text{ free} & : & 4 \\ \text{YOUR\_NAME} & : & -1 \\ \text{MISPELLED} & : & 1 \\ \text{FROM\_FRIEND} & : & -3 \\ \dots & & \end{pmatrix} w$$

$f(x_1)$

$$\begin{pmatrix} \# \text{ free} & : & 2 \\ \text{YOUR\_NAME} & : & 0 \\ \text{MISPELLED} & : & 2 \\ \text{FROM\_FRIEND} & : & 0 \\ \dots & & \end{pmatrix}$$

$f(x_2)$

$$\begin{pmatrix} \# \text{ free} & : & 0 \\ \text{YOUR\_NAME} & : & 1 \\ \text{MISPELLED} & : & 1 \\ \text{FROM\_FRIEND} & : & 1 \\ \dots & & \end{pmatrix}$$

- If correct (i.e., predicted  $y$ =target  $y^*$ ), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if  $y^*$  is -1.

$$w = w + y^* \cdot f$$

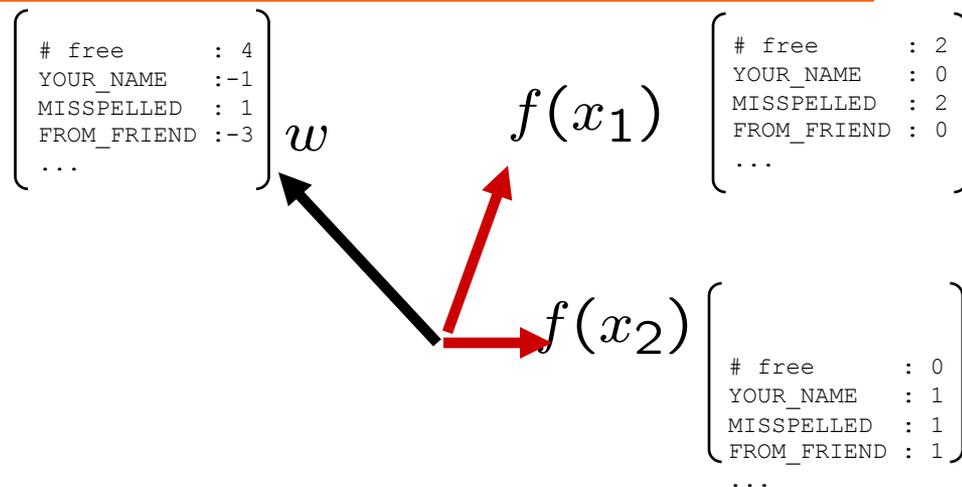
Why?

# SGD training for Binary Classifier

Figure out the weight vector from training instances

- Start with weights = 0
- For each training instance:
  - Classify with current weights
  - $f(x)$  is feature vector of  $x$

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$



## SGD with learning rate 1:

Do until satisfied:

- For each training example  $(y^*, f)$

1. Compute the gradient  $\nabla E$  where  $E$  is squared error

2. Update  $w = w - \nabla E$

*Namely no change with correct prediction*

*Otherwise  $w = w + y^* \cdot f$*

$$E = 0.5 (y^* - f(x))^2$$

$$\nabla E = \partial E / \partial w = -(y^* - y) f$$

$$= 0 \text{ if } y^* = y$$

$$\text{else } -y^* f$$

# Example of SGD Learning from training data

- *Classifier model:*  $f(x) = \text{Size} * w_1 + \text{color} * w_2 + \text{shape} * w_3$   
Use sign of  $f(x)$  to classify  
Initially  $w_1 = w_2 = w_3 = 0$

Instance	Size	Color	Shape	Category
$x_1$	Small 0	Red 0	Circle 0	Positive 1
$x_2$	Large 2	Red 0	Circle 0	Positive 1
$x_3$	Small 0	Red 0	Triangle 1	Negative -1
$x_4$	Large 2	Blue 1	Circle 0	Negative -1

With Instance 1:  $\text{sign}(f(x_1)) = \text{sign}(0) = 1$ . No weight change

With Instance 2:  $\text{sign}(f(x_2)) = \text{sign}(0) = 1$ . No weight change.

With Instance 3:  $\text{sign}(f(x_3)) = \text{sign}(0) = 1$ . Wrongly classified  
 $w = w + (-1) * (0, 0, 1) = (0, 0, -1)$

With Instance 4:  $\text{sign}(f(x_4)) = \text{sign}(0) = 1$ . Wrongly classified  
 $w = w + (-1) * (2, 1, 0) = (-2, -1, -1)$

# Optimization Problem for Classification

---

Given training set  $\{(x_1, y_1), \dots, (x_n, y_n)\}$

Given a loss function  $\ell(h, y)$  (hinge loss, logistic, ...)

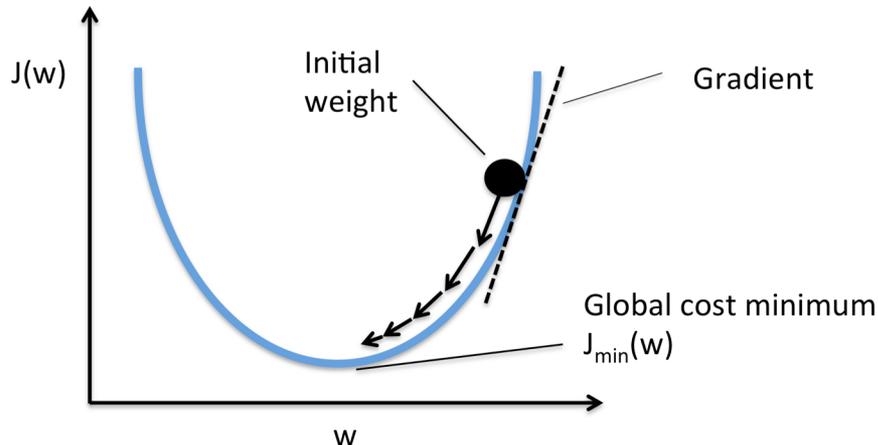
Find a prediction function  $h(x; w)$  (linear, DNN, ...)

$$\min_w \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w), y_i)$$

- "y<sub>i</sub>" is the classification label for a training instance
- "w" is the set of parameters to be found through training
- What does prediction function h() look like?
- How to find parameters involved in h() that minimize an objective function?

# How to find parameters that minimize the loss function?

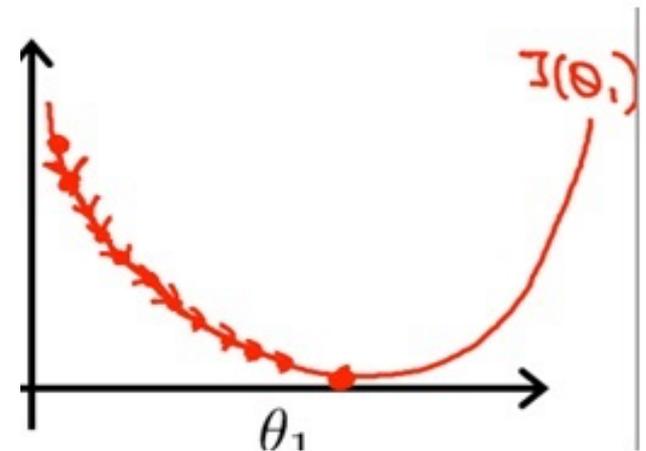
- How to find parameters that minimize a loss function  $J$  with parameter vector  $w^{\mathcal{D}}$



- Gradient Descent Method (SGD) for Optimization
  - Start somewhere
  - Repeat: Take a step in the steepest descent direction

Learning rate

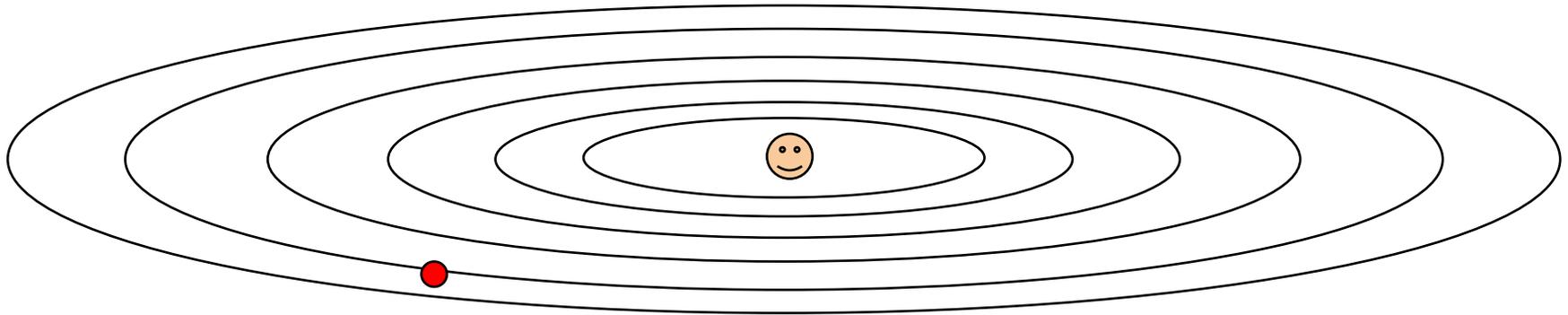
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$



$w$  is  $\theta_1$

# Illustration of gradient descent to refine multiple parameters

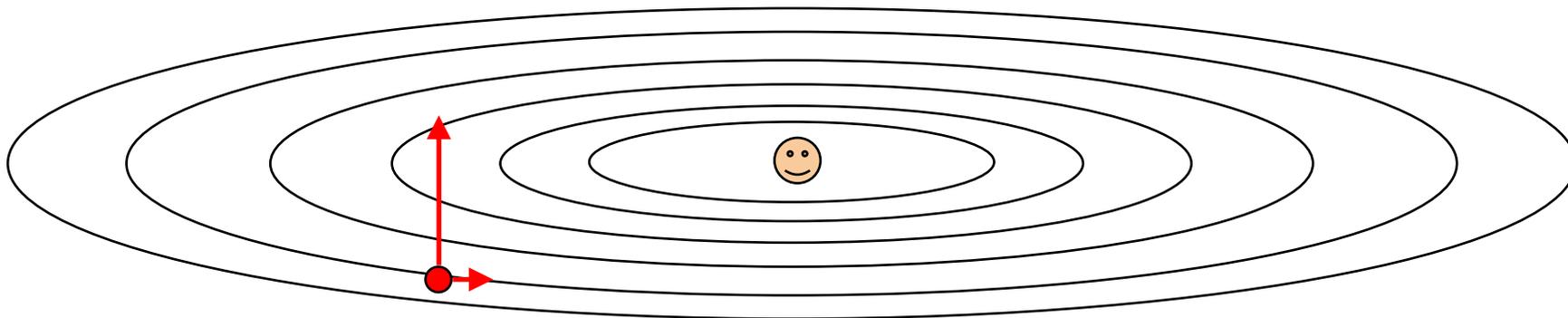
---



Q: What is the trajectory along which we converge towards the minimum with SGD?

Suppose loss function is steep vertically but shallow horizontally:

---

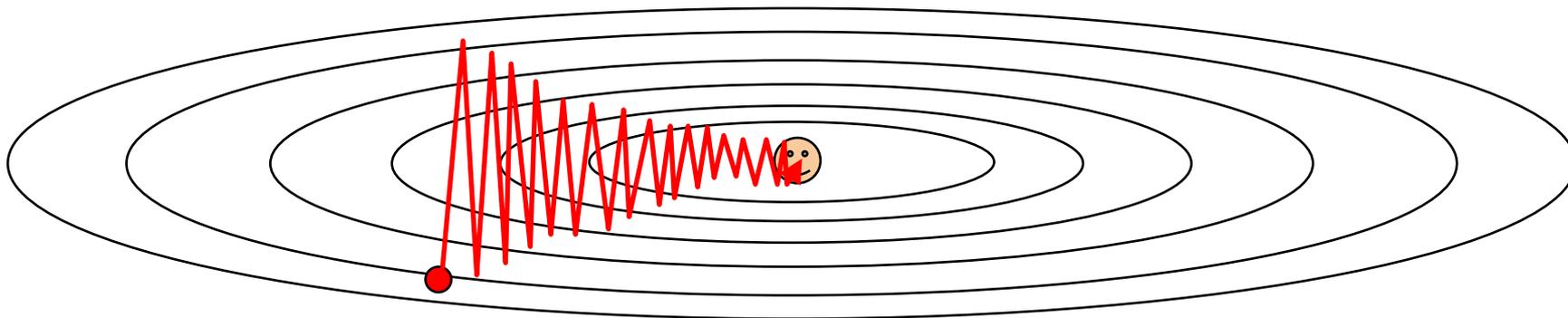


Q: What is the trajectory along which we converge towards the minimum with SGD?

10

Suppose loss function is steep vertically but shallow horizontally:

---



Q: What is the trajectory along which we converge towards the minimum with Gradient Descent? **very slow progress** along flat direction, **jitter** along steep one

11

# Generally, Steepest Direction with $n$ parameters

- Given loss function  $g$  and learning rate  $\alpha$

- Steepest Direction = direction of the gradient

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \dots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

- Parameter vector  $w = (w_1, w_2, \dots, w_n)$

- Gradient Descent: Update weight vector  $w$  by using a sequence of training instance  $i$

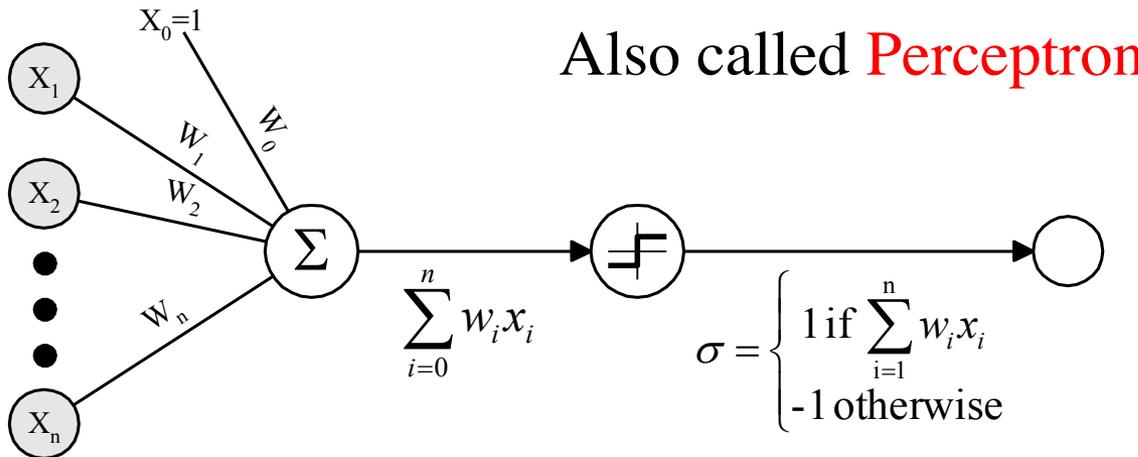
1. Stop after a fixed number of iterations.
2. Or when loss is close to a lower bound or has not improved much in a long time.
3. Or when the validation error has not improved in a long time.

- Init:

- For  $i = 1, 2, \dots$

$$w \leftarrow w - \alpha * \nabla g(w)$$

# Linear Classifier with Square Error



$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

Sometimes we will use simpler vector notation :

$$o(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise} \end{cases}$$

# SGD for Linear Classifier (Perceptron)

To understand, consider simple *linear unit*, where

$$o = w_0 + w_1 x_1 + \dots + w_n x_n$$

Idea : learn  $w_i$ 's that minimize the squared error

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

Where  $D$  is the set of training examples

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta (t - o) x_i$$

- $t = c(\vec{x})$  is target value
- $o$  is perceptron output
- $\eta$  is small constant (e.g., .1) called learning rate

Can prove it will converge

- If training data is linearly separable
- and  $\eta$  is sufficiently small

Training:  $w_i = w_i - \partial E / \partial w_i$

What is  $\partial E / \partial w_i$  ?

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \end{aligned}$$

$$\frac{\partial E}{\partial w_i} = \sum_d (t_d - o_d) (-x_{i,d})$$

# Incremental vs Batch Mode in SGD

---

## SGD in an incremental mode:

Update weights instance by instance

Do until satisfied:

- For each training example  $d$  in  $D$

1. Compute the gradient  $\nabla E_d[\vec{w}]$

2.  $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d[\vec{w}]$

$$E_d[\vec{w}] \equiv \frac{1}{2} (t_d - o_d)^2$$

$$\nabla E = \partial E / \partial \omega = -(t_d - o_d)x$$

---

## SGD in a batch or minibatch mode:

Update weights by a (mini-) batch of instances (subset  $D$ )

Do until satisfied:

1. Compute the gradient  $\nabla E_D[\vec{w}]$

2.  $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D[\vec{w}]$

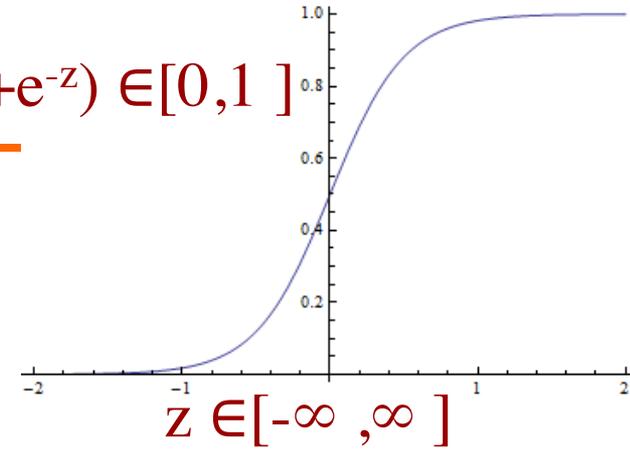
$$E_D[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

---

*Batch Gradient is a good approximation with small learning rate while allowing faster computation*

# Other Classification Prediction or Loss Functions

$$e^z / (e^z + e^{-z}) \in [0, 1]$$



## Softmax for binary classification

## Logistic regression

- Score for  $y=1$ :  $w^\top f(x)$

- Score for  $y=-1$ :  $-w^\top f(x)$

- Probability of label:

$$p(y = 1 | f(x); w) = \frac{e^{w^\top f(x^{(i)})}}{e^{w^\top f(x)} + e^{-w^\top f(x)}}$$

$$p(y = -1 | f(x); w) = \frac{e^{-w^\top f(x)}}{e^{w^\top f(x)} + e^{-w^\top f(x)}}$$

- Maximize: 
$$l(w) = \prod_{i=1}^m p(y = y^{(i)} | f(x^{(i)}); w)$$

- Equivalently maximize log likelihood:

$$ll(w) = \sum_{i=1}^m \log p(y = y^{(i)} | f(x^{(i)}); w)$$

# Multi-class Softmax

---

- 3-class softmax – classes A, B, C
  - 3 weight vectors:  $w_A, w_B, w_C$

- Probability of label A: (similar for B, C)

$$p(y = A | f(x); w) = \frac{e^{w_A^\top f(x)}}{e^{w_A^\top f(x)} + e^{w_B^\top f(x)} + e^{w_C^\top f(x)}}$$

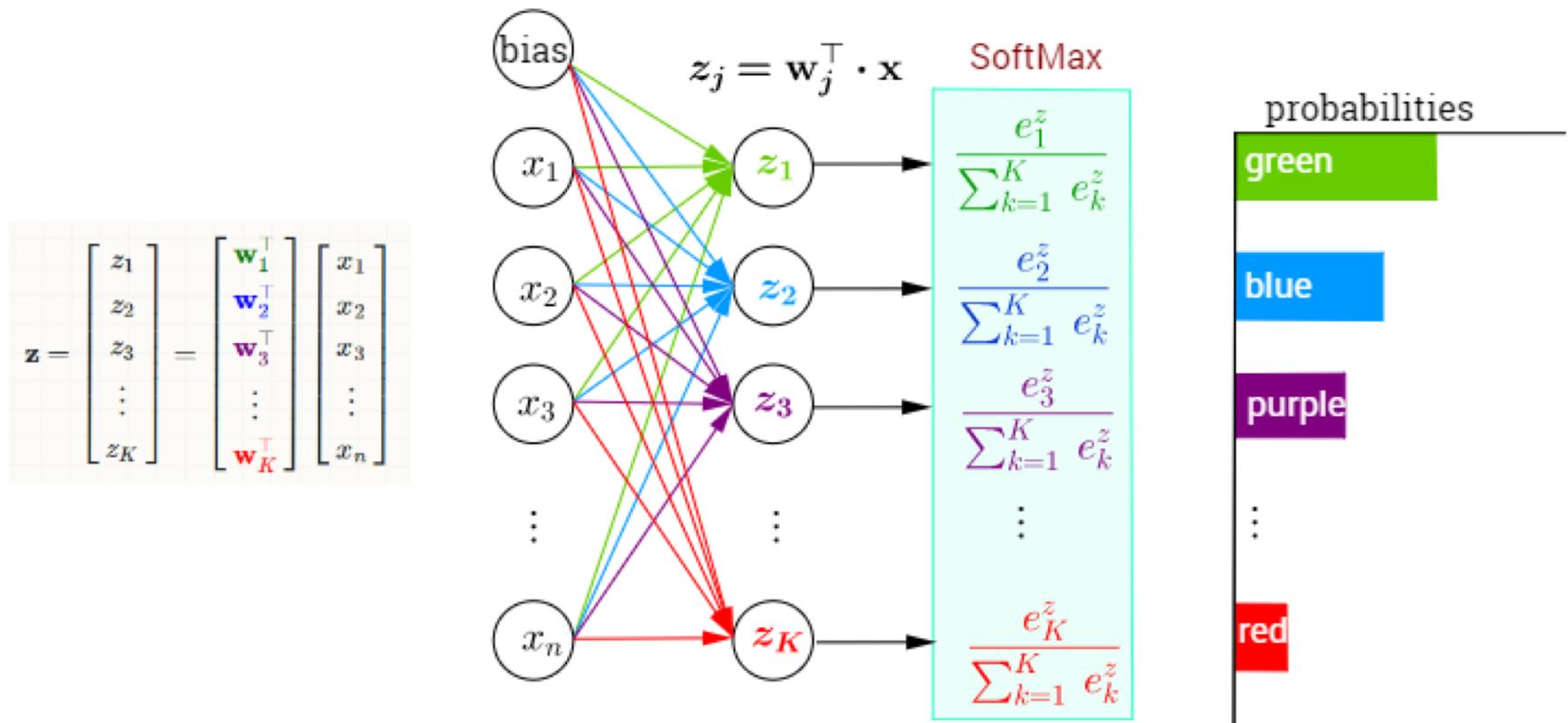
- Loss function:  $l(w) = \prod_{i=1}^m p(y = y^{(i)} | f(x^{(i)}; w))$

- Equivalently maximize log likelihood:

$$ll(w) = \sum_{i=1}^m \log p(y = y^{(i)} | f(x^{(i)}; w))$$

# Multi-class Two-Layer Neural Network with SoftMax

## Multi-Class Classification with NN and SoftMax Function



# Activation Function: tanh(x)

## Graphical Representation of tanh(x)

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 1 - \frac{2e^{-x}}{e^x + e^{-x}}$$

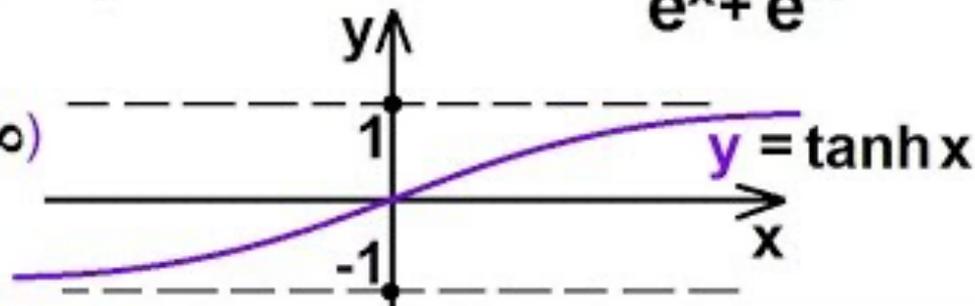
$$\text{as } x \rightarrow \infty \frac{2e^{-x}}{e^x + e^{-x}} \rightarrow 0$$

$$\text{when } x=0 \frac{2e^{-x}}{e^x + e^{-x}} = 1$$

$$\text{as } x \rightarrow -\infty \frac{2e^{-x}}{e^x + e^{-x}} \rightarrow 2$$

$$e^x + e^{-x} \sqrt{e^x - e^{-x}}$$

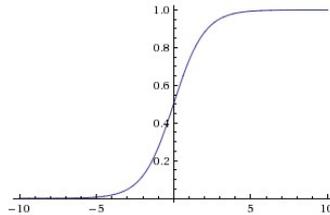
odd function  
domain =  $(-\infty, \infty)$   
range =  $(-1, 1)$



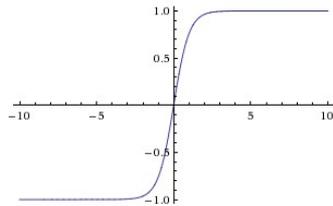
# Other Activation Functions

## Sigmoid

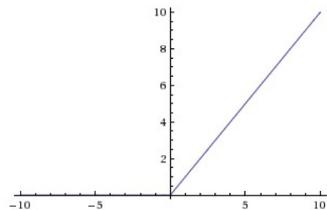
$$\sigma(x) = 1/(1 + e^{-x})$$



## tanh tanh(x)

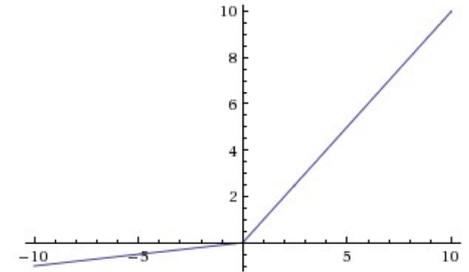


## ReLU max(0,x)



## Leaky ReLU

$$\max(0.1x, x)$$

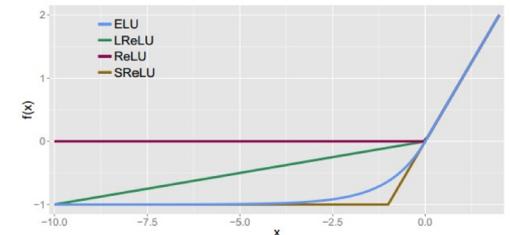


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

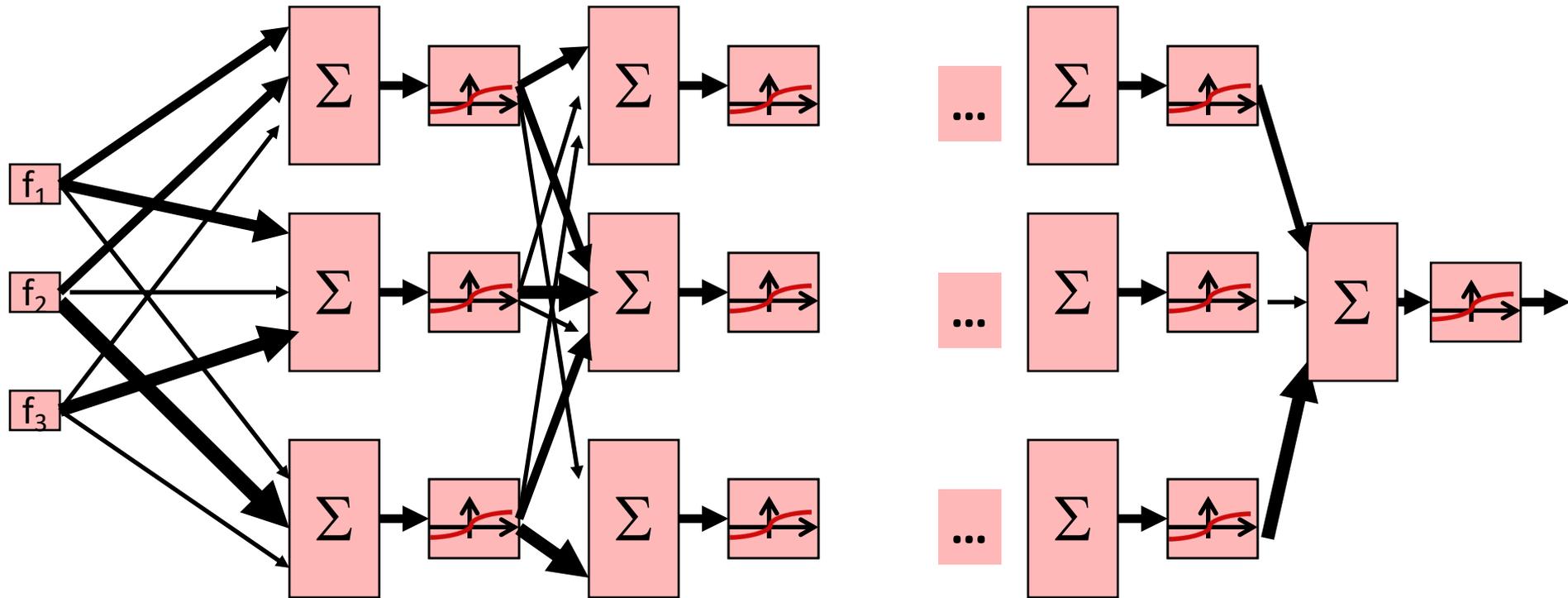
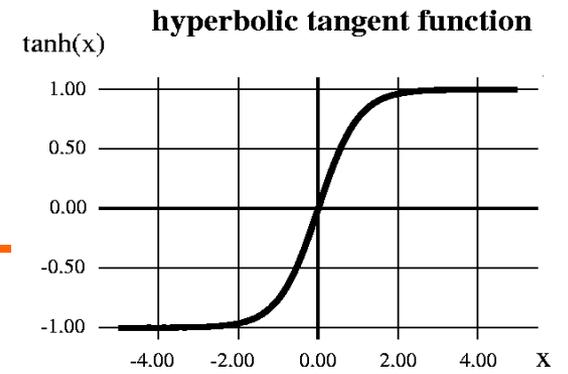
## ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



20

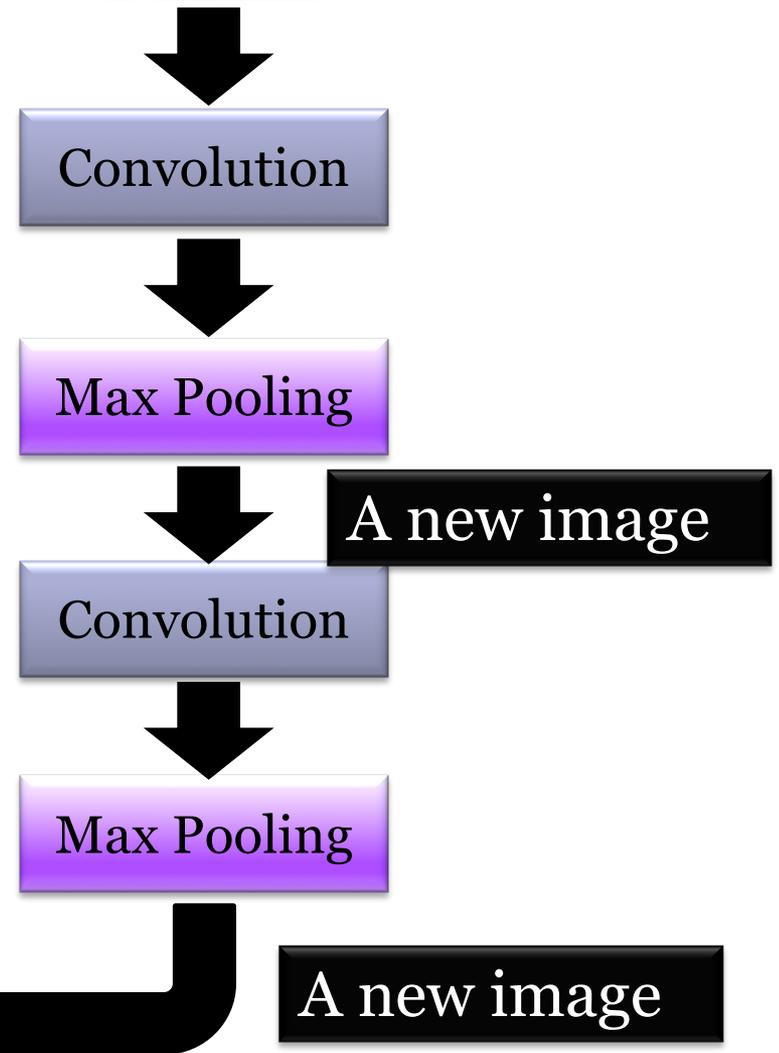
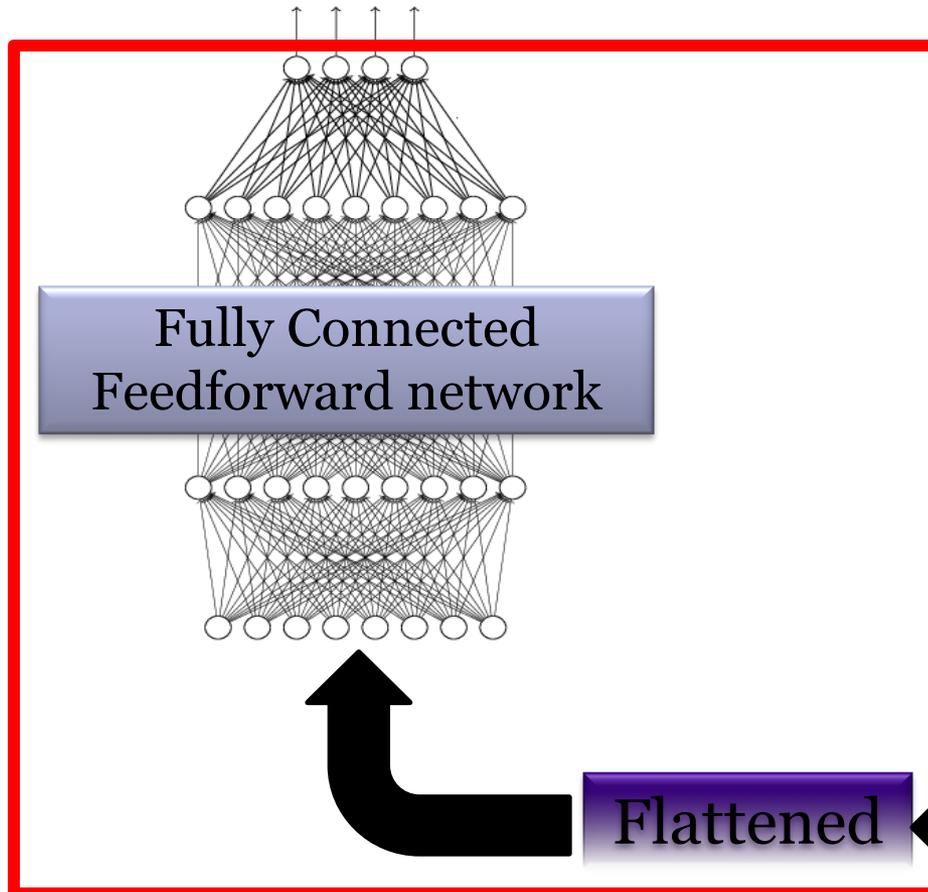
# N-Layer Neural Network



# The whole CNN



cat dog .....



# How to Calculate Partial Derivatives for SGD through a Computer Algorithm

---

- Graph representation of a loss function can be huge with thousands or even millions of parameters.
- How to compute partial derivatives of a computational graph

Example: Given a function  $f(x,y,z) = (x+y)z$ , what is the partial derivative of  $f$  with respect to  $x$ ,  $y$ ,  $z$ ?

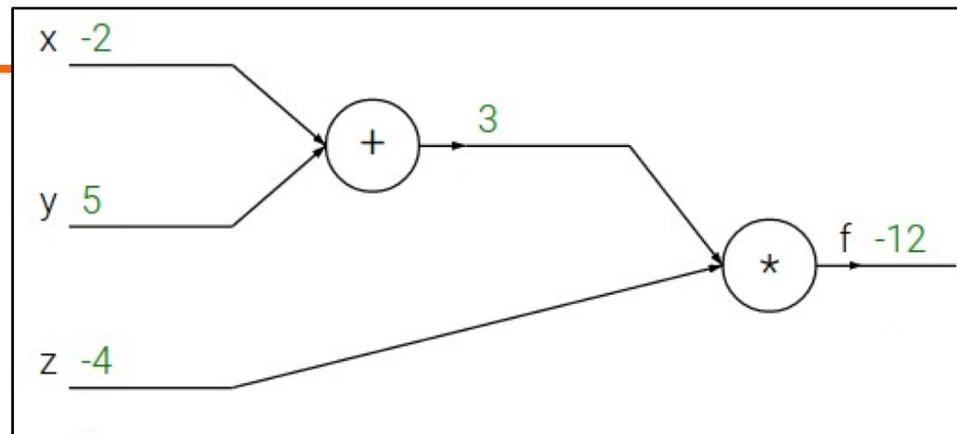
- Computer has to do it symbolically. Not easy in general
- What is the partial derivative of  $f$  with respect to  $x$ ,  $y$ ,  $z$ , given  $x = -2$ ,  $y = 5$ ,  $z = -4$  from a training instance?

**Easier to do by focusing on the given training instance**

# Example of Algorithmic Derivative Computation

-  $f(x, y, z) = (x + y)z$

Knowing  $x = -2, y = 5, z = -4$

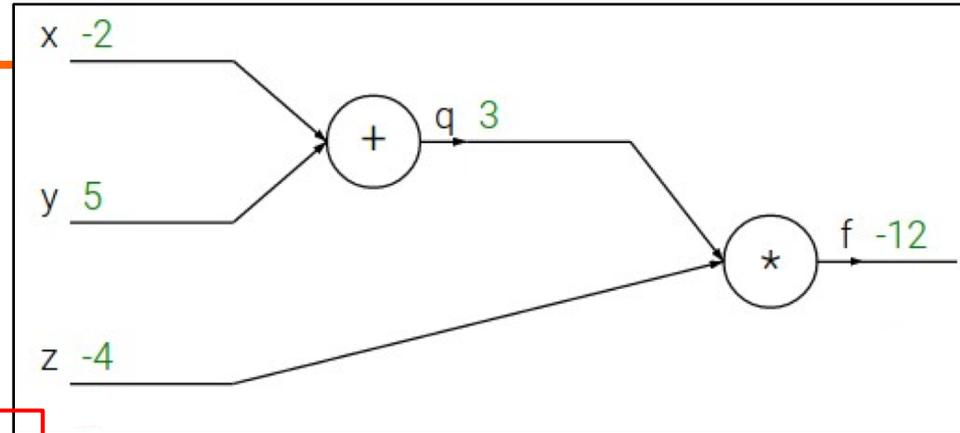


# Get local derivatives for each node

# Get the final value $f$ via forward computation

-  $f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$



Get local derivatives for each node

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

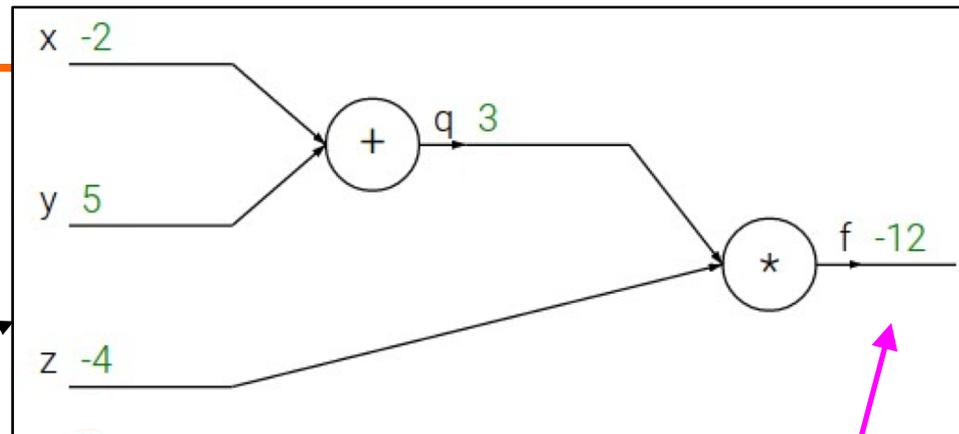
Now we conduct a backward propagation in this graph to compute  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Backward to get the derivative of last node

$$\frac{\partial f}{\partial f}$$

•  $f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$



$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

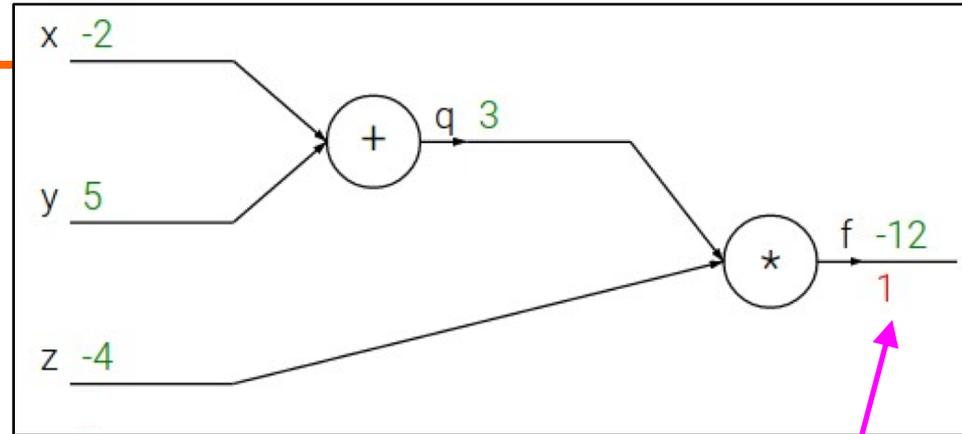
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

$$\frac{\partial f}{\partial f}$$

$\frac{\partial f}{\partial f} = 1$  as local derivative. It is trivial

$f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$



$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$

$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Need to get derivative

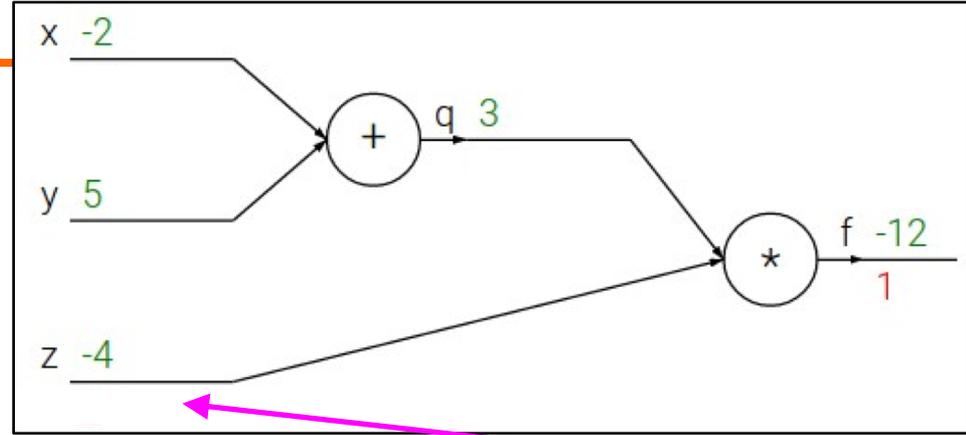
$$\frac{\partial f}{\partial z}$$

$f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x,y,z) = -12$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$



$$\frac{\partial f}{\partial z}$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Derive 3 as derivative

$$\frac{\partial f}{\partial z}$$

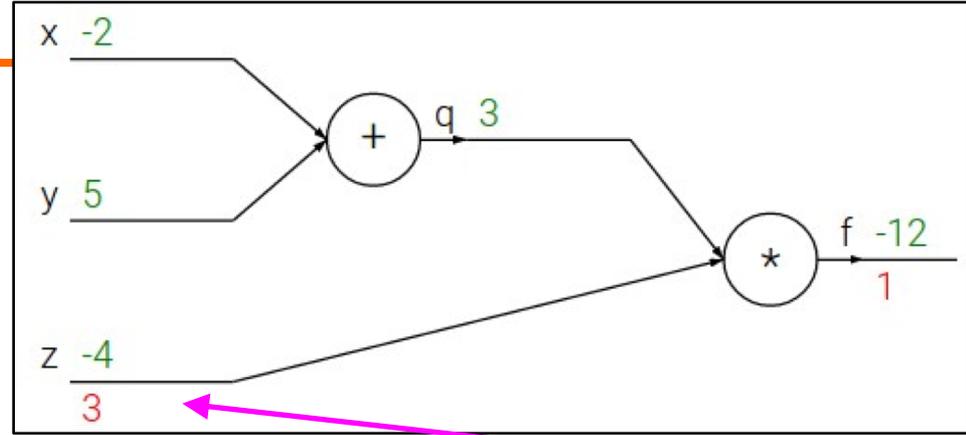
because  $\partial f / \partial z = q = 3$

$f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$



$$\frac{\partial f}{\partial z}$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Need to get derivative

$$\frac{\partial f}{\partial q}$$

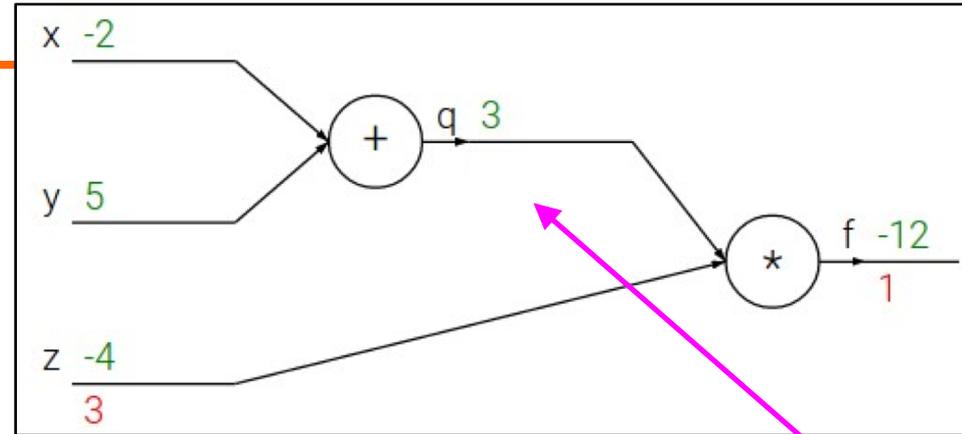
$$f(x, y, z) = (x + y)z$$

$$x = -2, y = 5, z = -4, f(x, y, z) = -12$$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

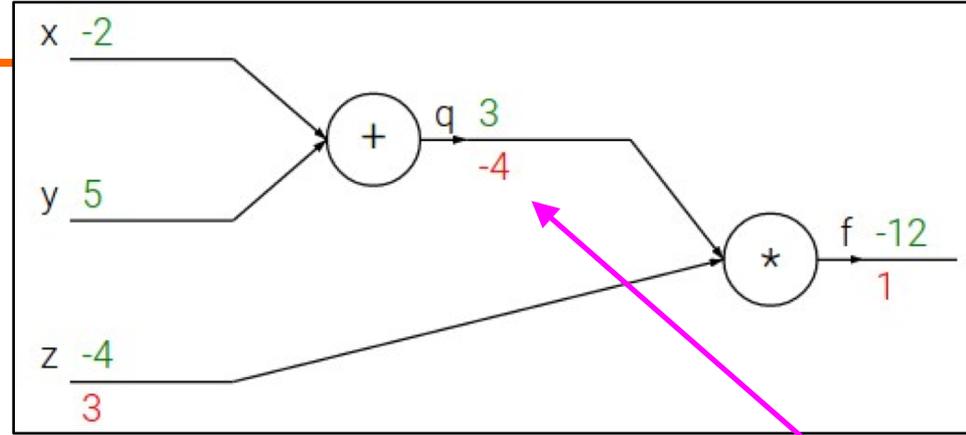
$\frac{\partial f}{\partial q}$  is found because  $\partial f / \partial q = z = -4$

$f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$

$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$

$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$



$\frac{\partial f}{\partial q}$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# How to compute $\partial f / \partial y$ ?

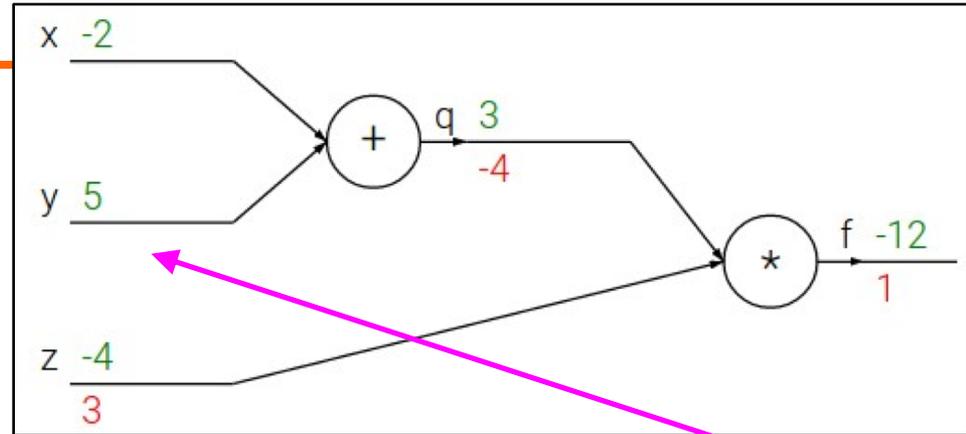
•  $f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

# Use the chain rule locally to compute $\frac{\partial f}{\partial y} = (-4) \cdot 1 = -4$

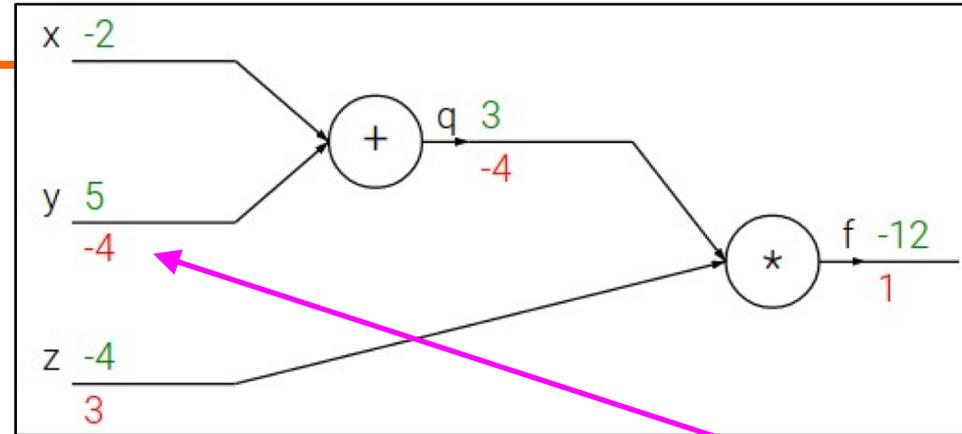
$f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$

$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$

$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$\frac{\partial f}{\partial y}$

Chain rule:

$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$

# Use the chain rule locally to compute $\partial f / \partial x$

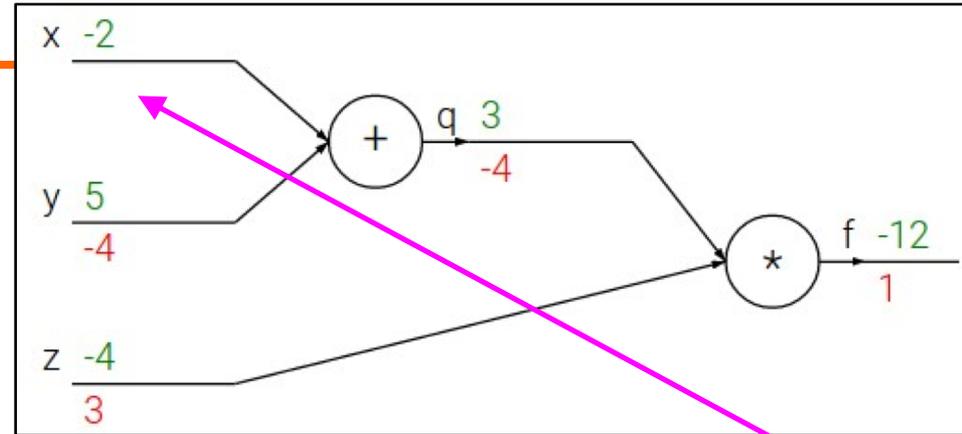
$f(x, y, z) = (x + y)z$

$x = -2, y = 5, z = -4, f(x, y, z) = -12$

$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$

$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$\frac{\partial f}{\partial x}$

# Use the chain rule locally to compute $\partial f / \partial x = (-4) \cdot 1 = -4$

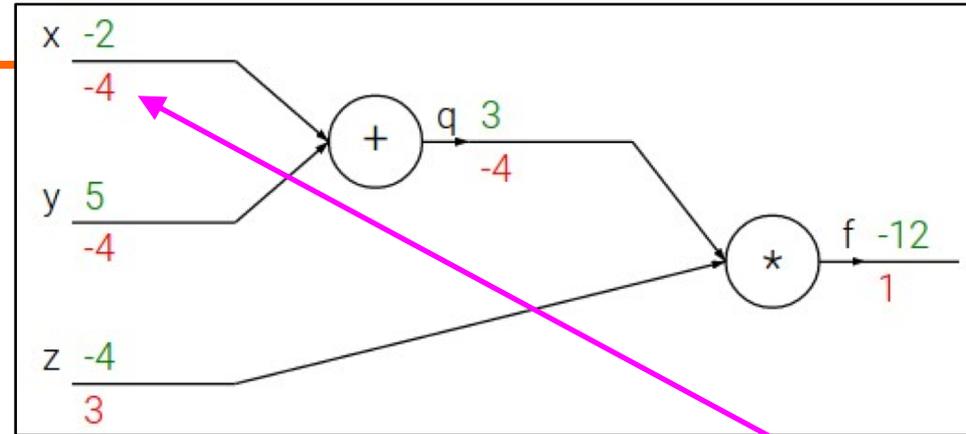
$$f(x, y, z) = (x + y)z$$

$$x = -2, y = 5, z = -4, f(x, y, z) = -12$$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

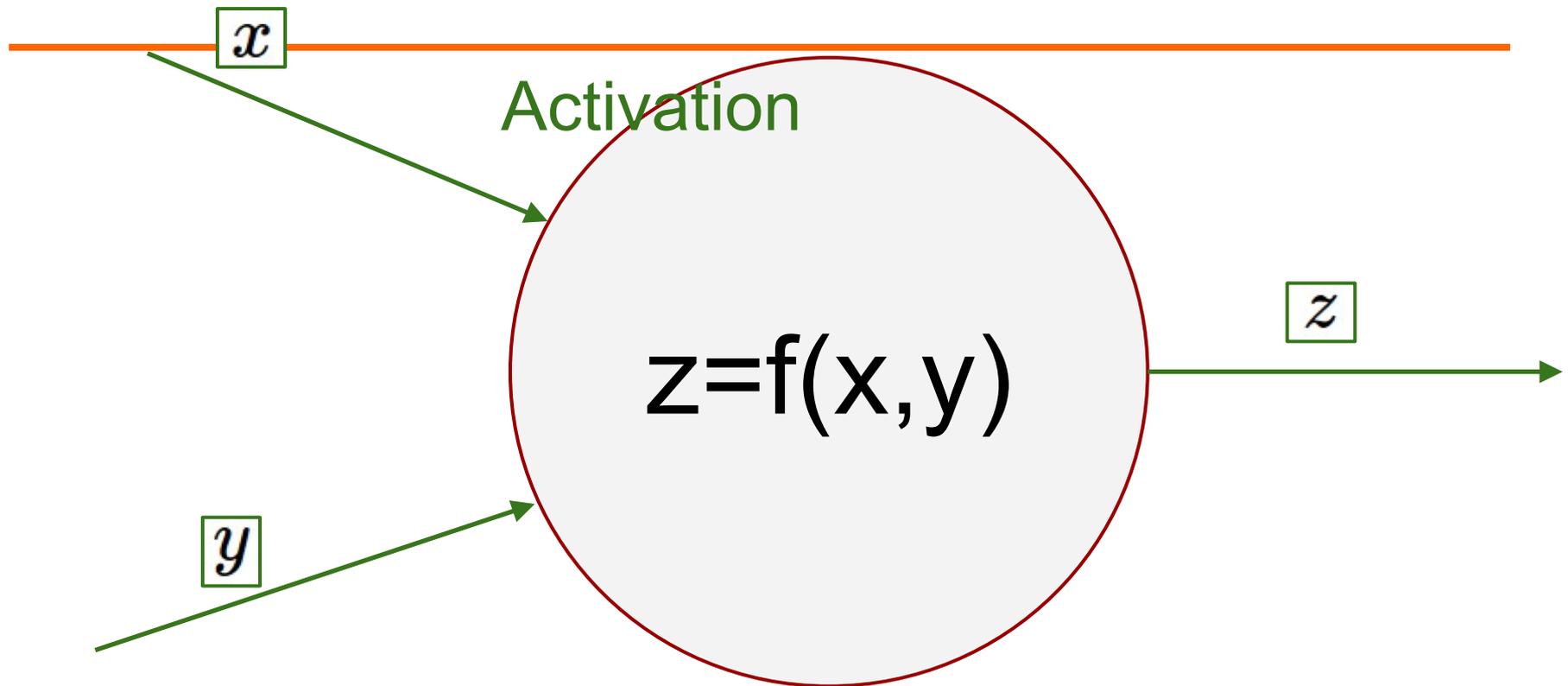


$$\frac{\partial f}{\partial x}$$

Chain rule:

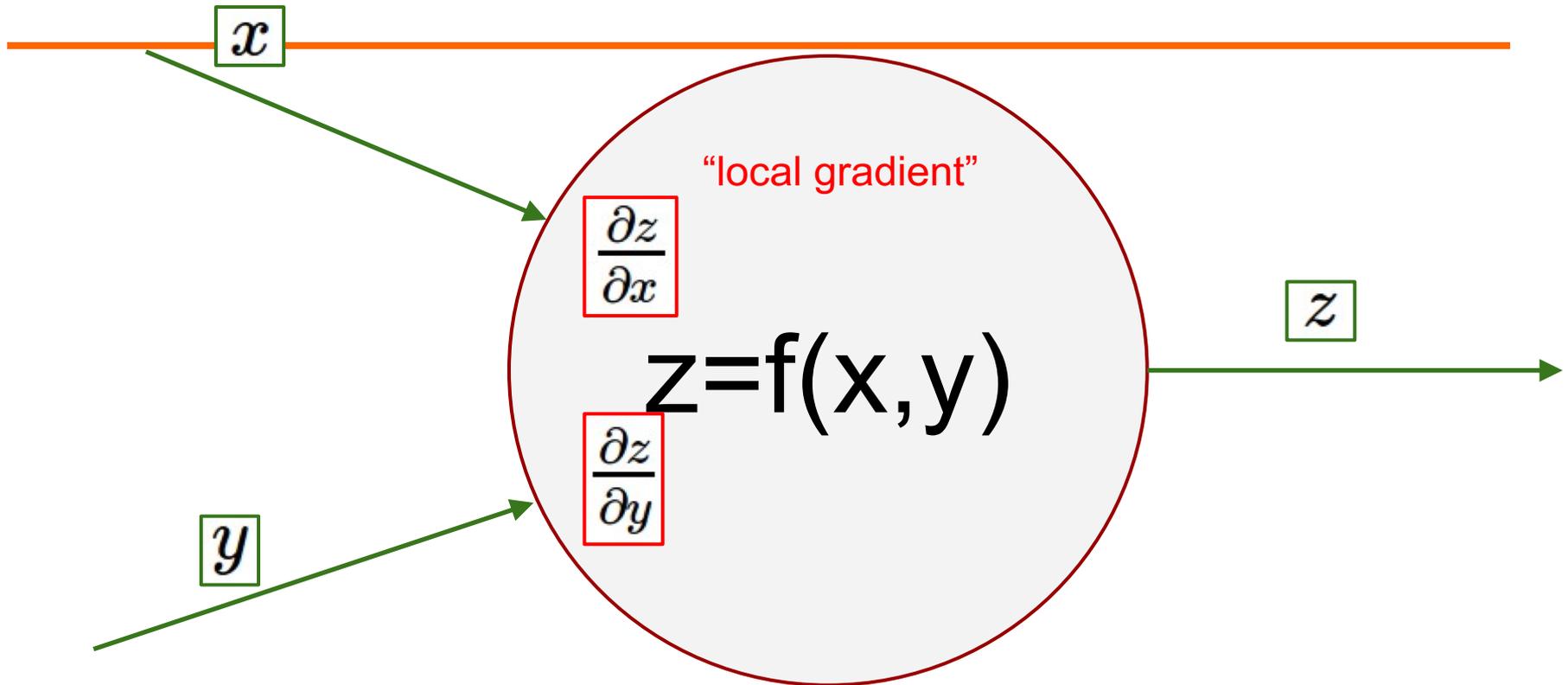
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

# How to use the chain rule locally ?



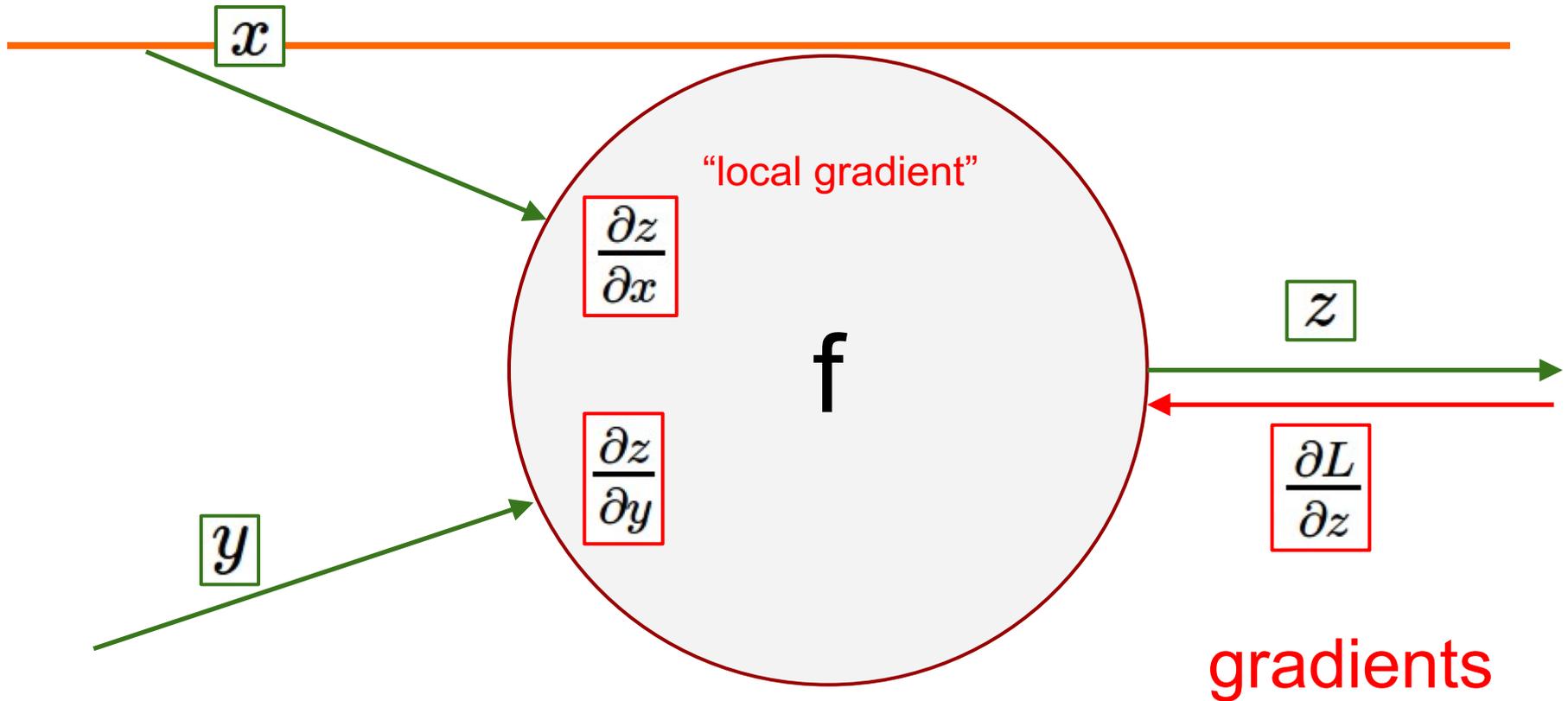
36

# Compute the local gradients first



37

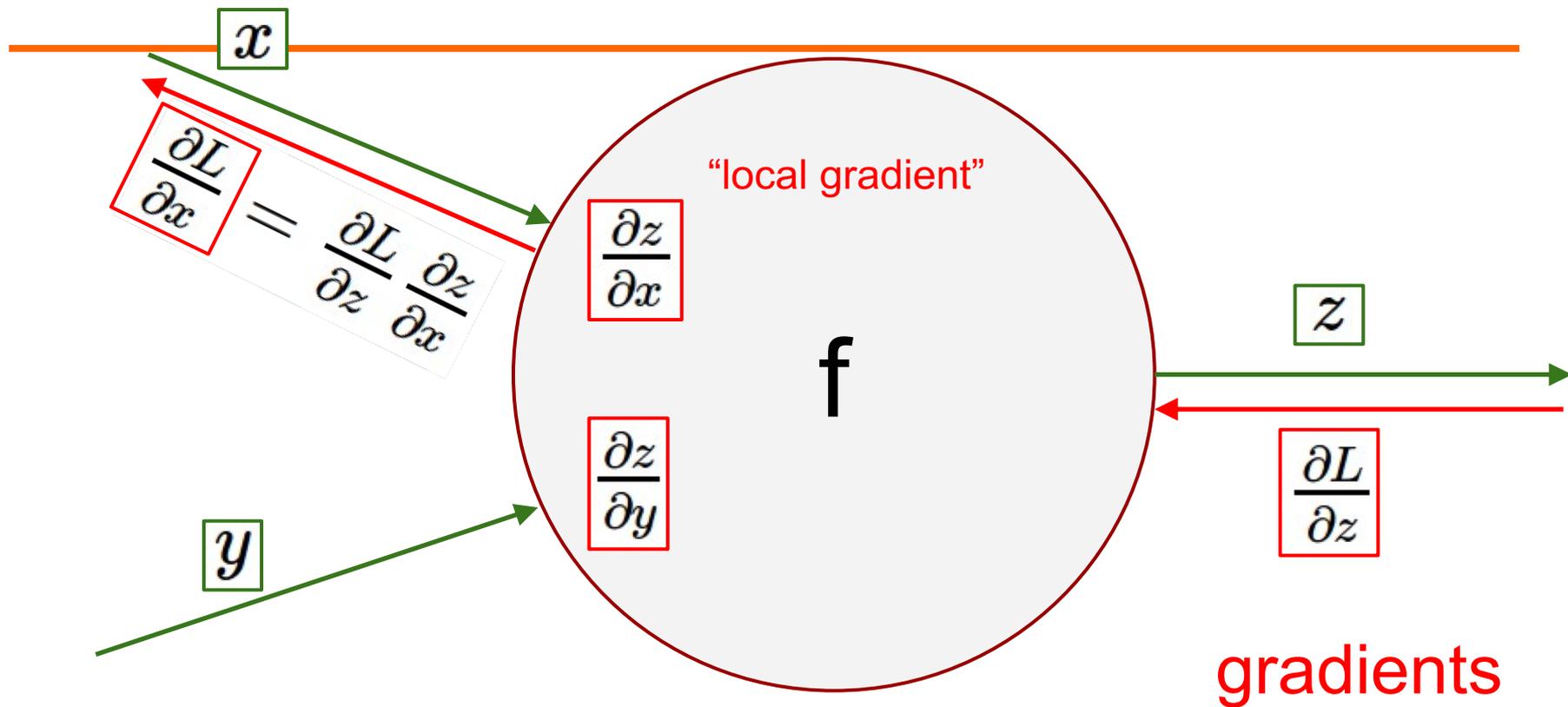
# Get the incoming gradient



38

# Apply the chain rule to compute the gradient

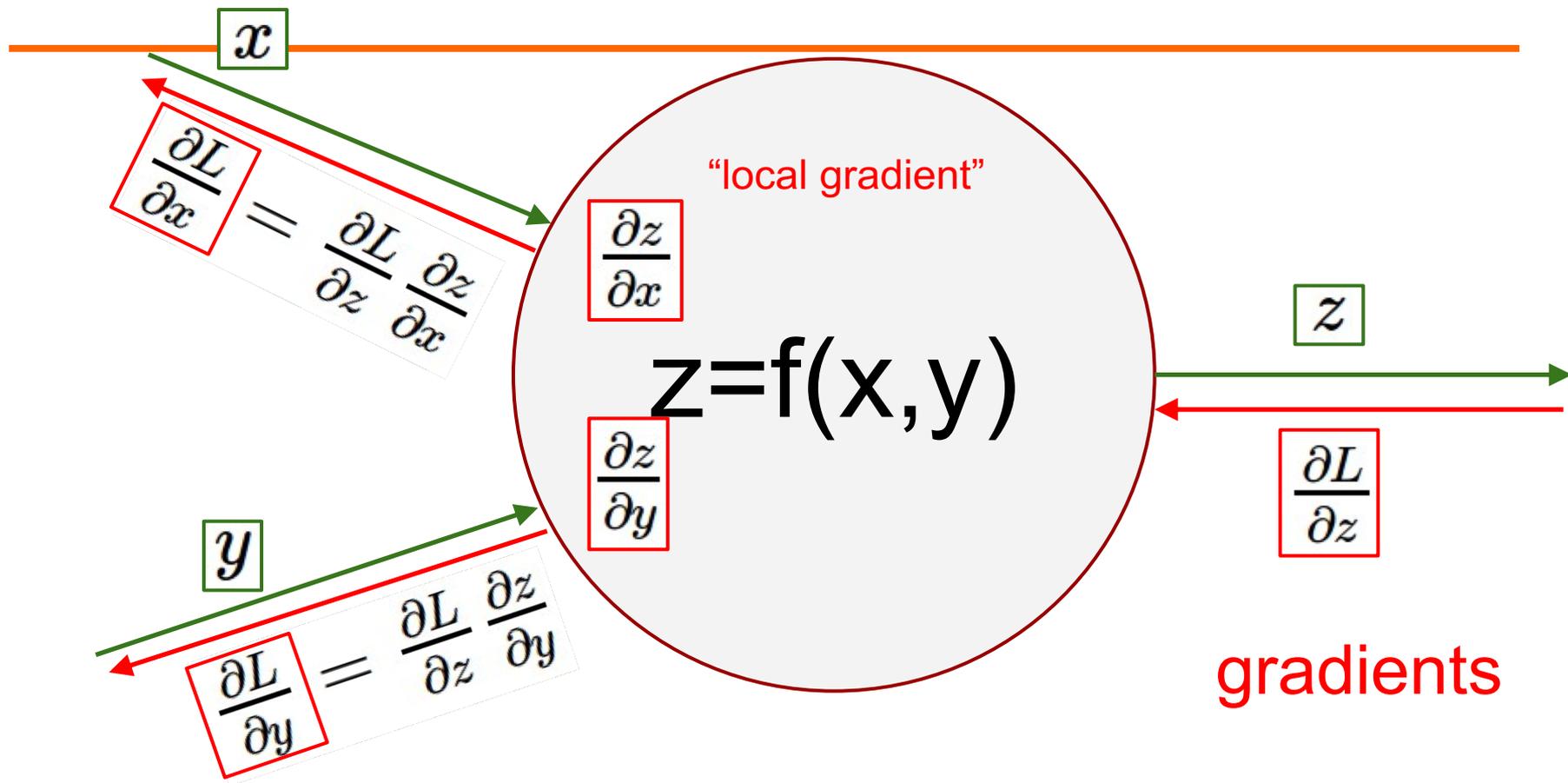
## Then propagate backward to one direction



39

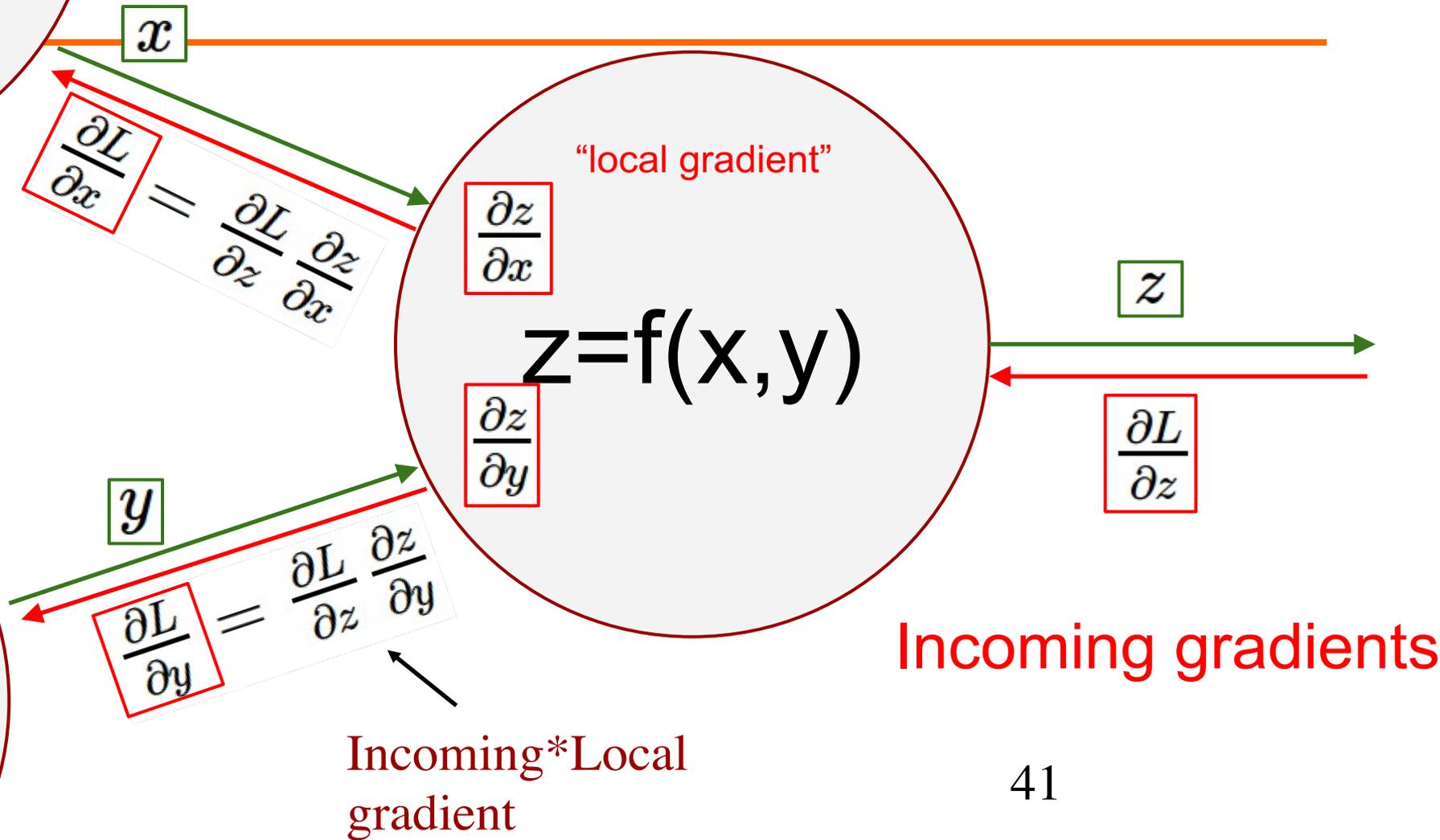
# Apply the chain rule to compute the gradient

## Propagate backward to another direction



40

# Summary of backward flow



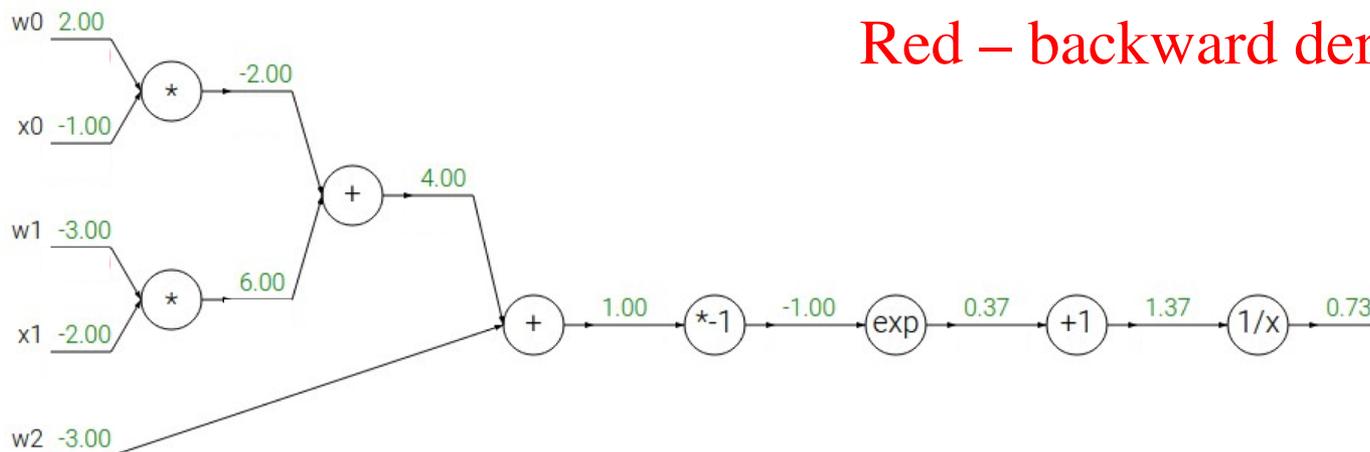
41

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Green- forward computation

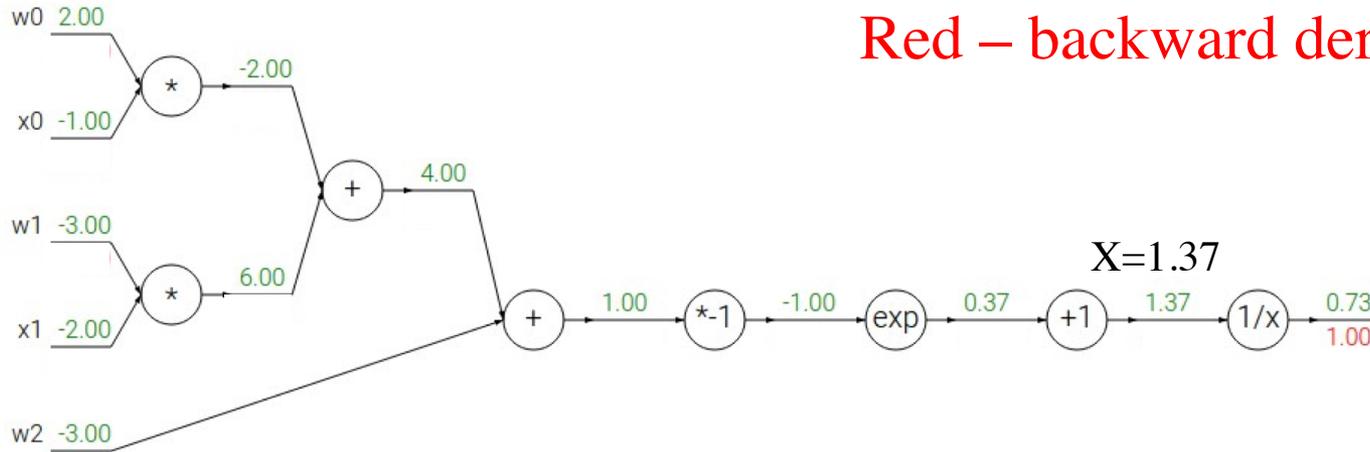
Red – backward derivatives



Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation  
Red – backward derivatives

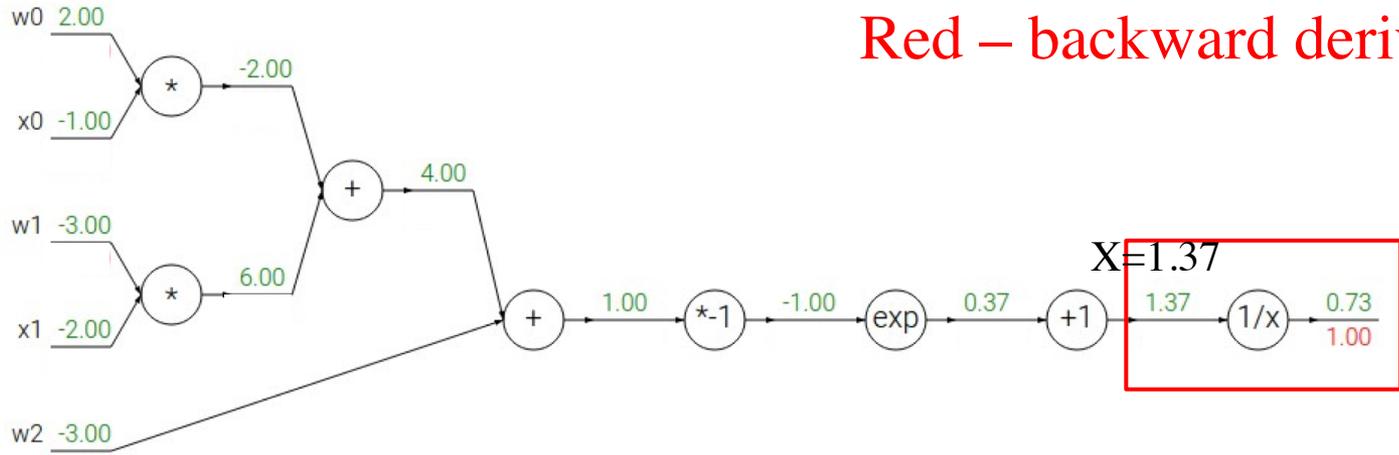


$f(x) = e^x$	→	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	→	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	→	$\frac{df}{dx} = a$		$f_c(x) = c + x$	→	$\frac{df}{dx} = 1$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation  
Red – backward derivatives



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

$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

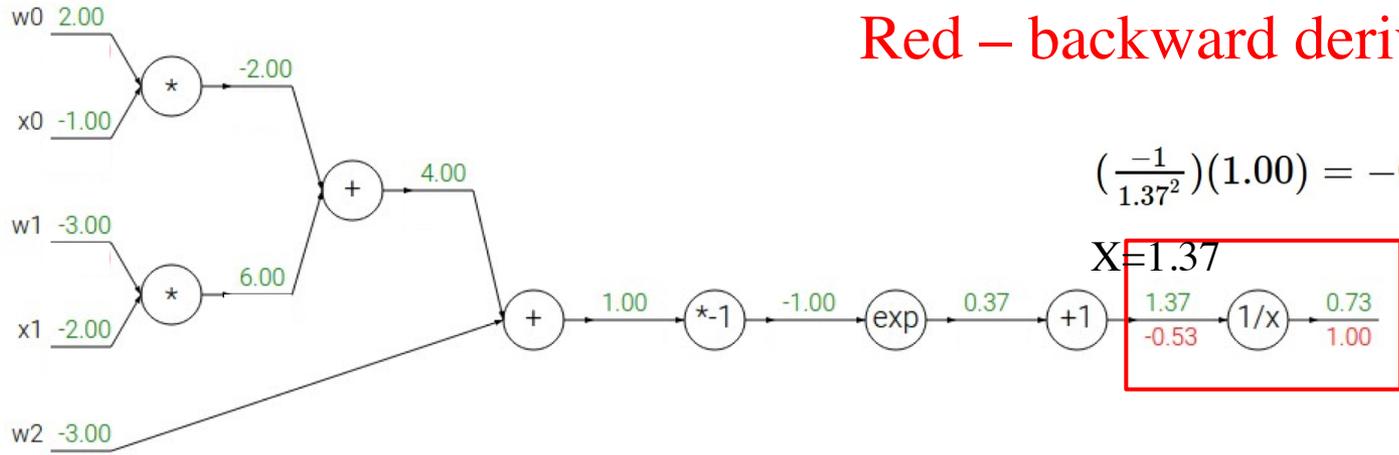
$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation  
Red – backward derivatives



$$\left(\frac{-1}{1.37^2}\right)(1.00) = -0.53$$

X=1.37

1.37	0.73
-0.53	1.00

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

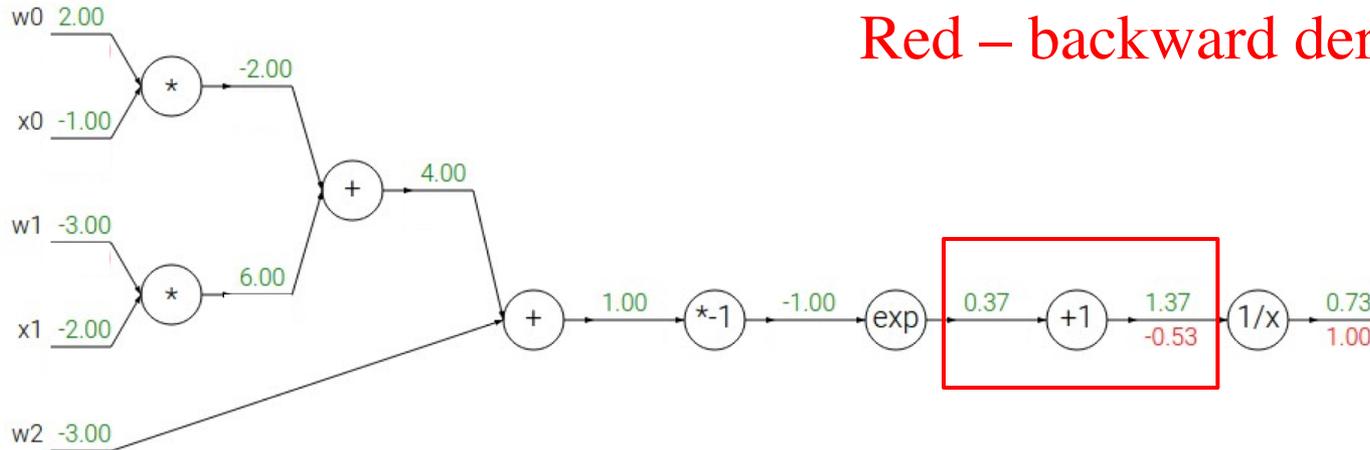
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$
$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation  
Red – backward derivatives



$f(x) = e^x$	$\rightarrow$	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	$\rightarrow$	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	$\rightarrow$	$\frac{df}{dx} = a$		$f_c(x) = c + x$	$\rightarrow$	$\frac{df}{dx} = 1$

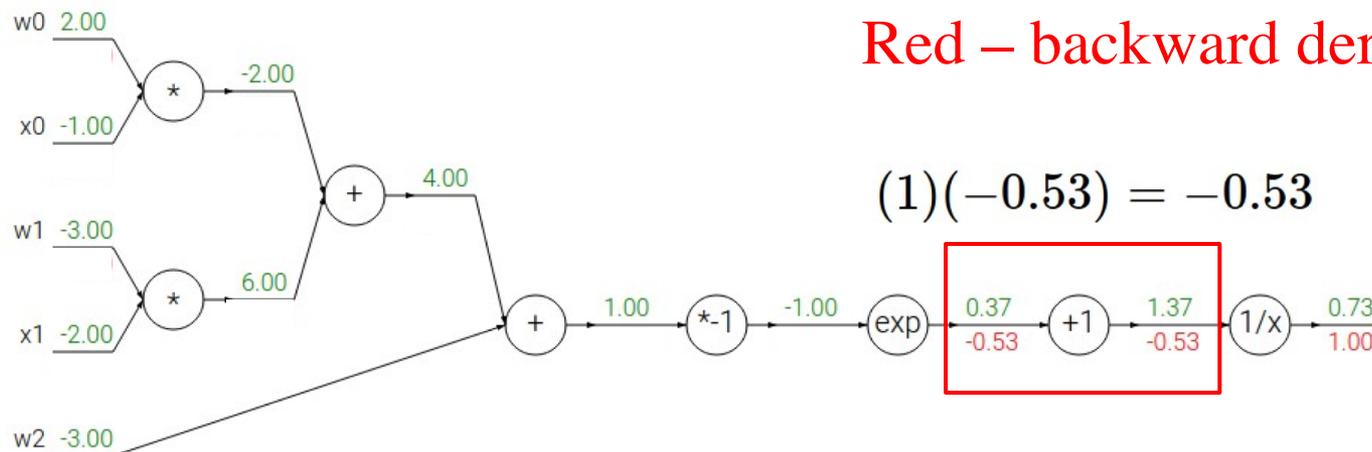
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation

Red – backward derivatives



$$(1)(-0.53) = -0.53$$

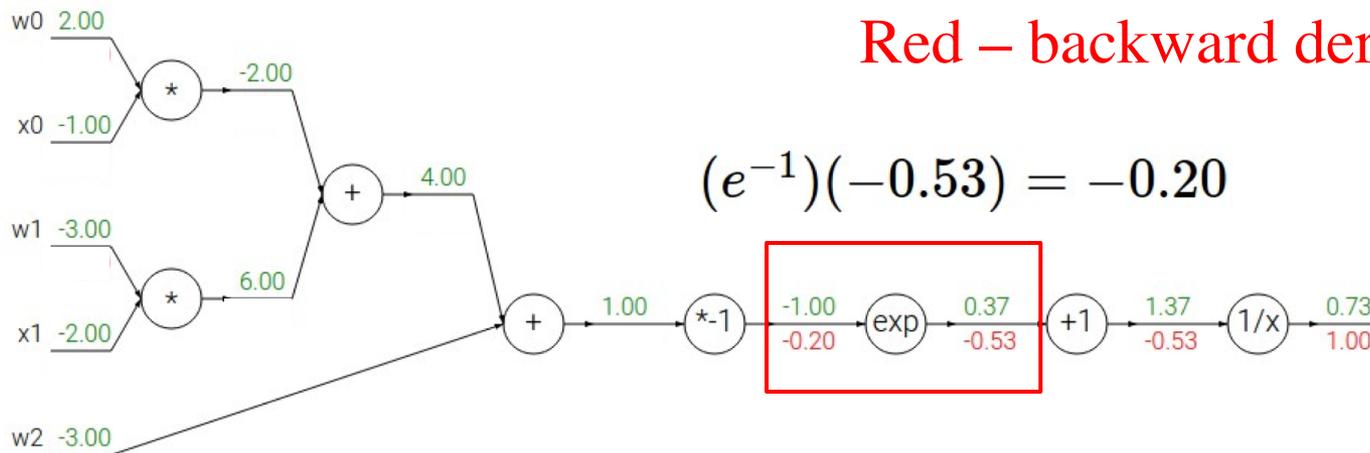
$f(x) = e^x$	→	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	→	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	→	$\frac{df}{dx} = a$		$f_c(x) = c + x$	→	$\frac{df}{dx} = 1$



Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation  
Red – backward derivatives



$$(e^{-1})(-0.53) = -0.20$$

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

$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

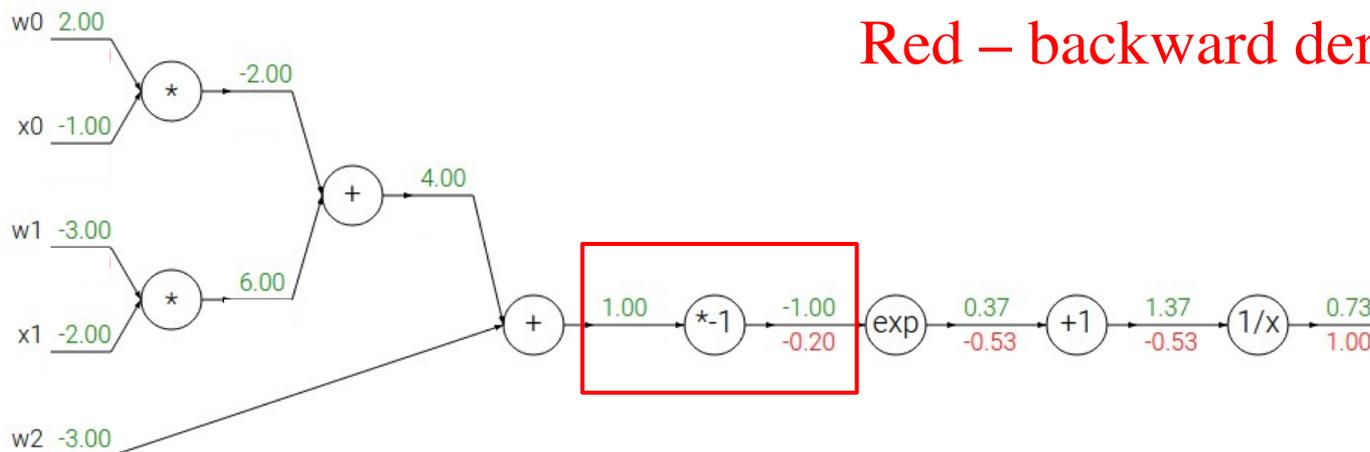
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation

Red – backward derivatives



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

$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

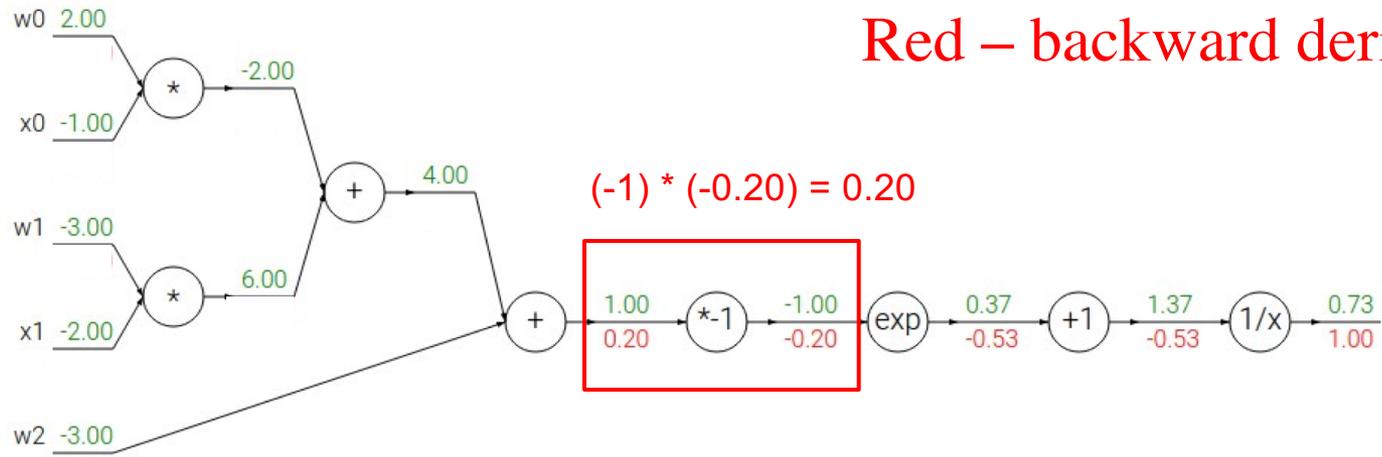
50

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green- forward computation  
Red – backward derivatives



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

$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

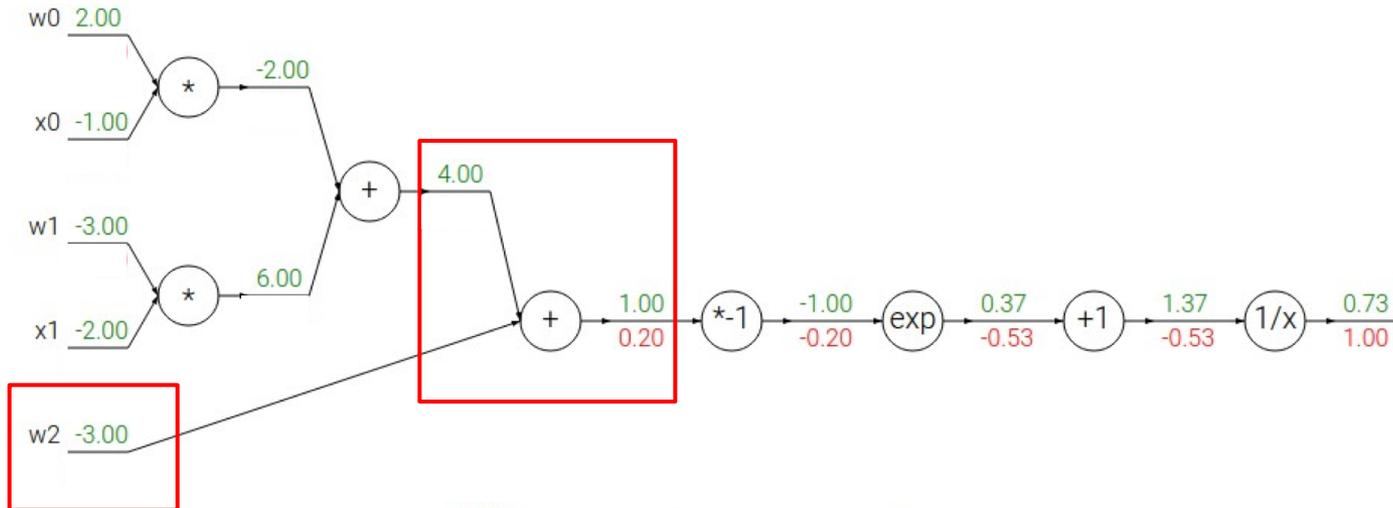
$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$$



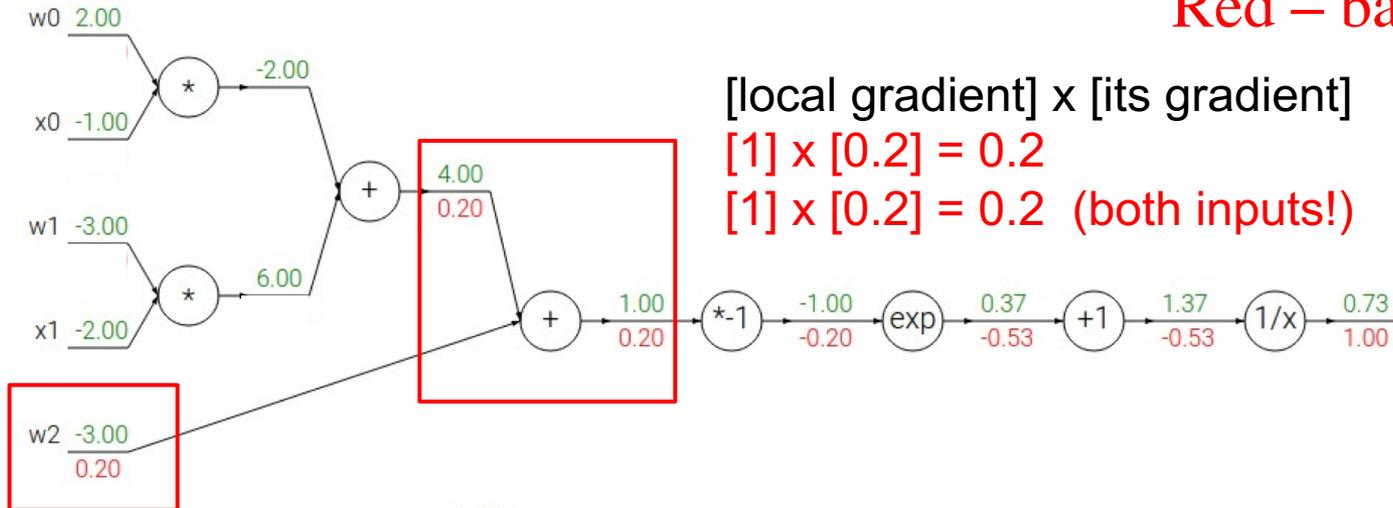
$f(x) = e^x$	→	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	→	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	→	$\frac{df}{dx} = a$		$f_c(x) = c + x$	→	$\frac{df}{dx} = 1$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Green - forward  
Red - backward



[local gradient] x [its gradient]  
[1] x [0.2] = 0.2  
[1] x [0.2] = 0.2 (both inputs!)

$$f(x) = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = e^x$$

$$\frac{df}{dx} = a$$



$$f(x) = \frac{1}{x}$$

$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

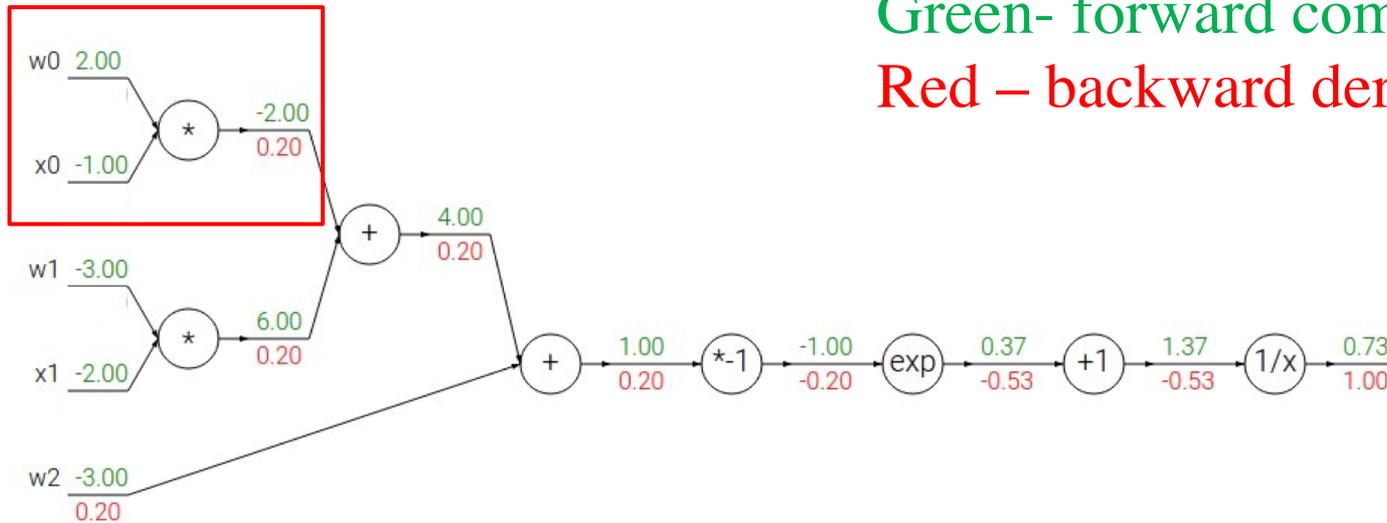
53

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

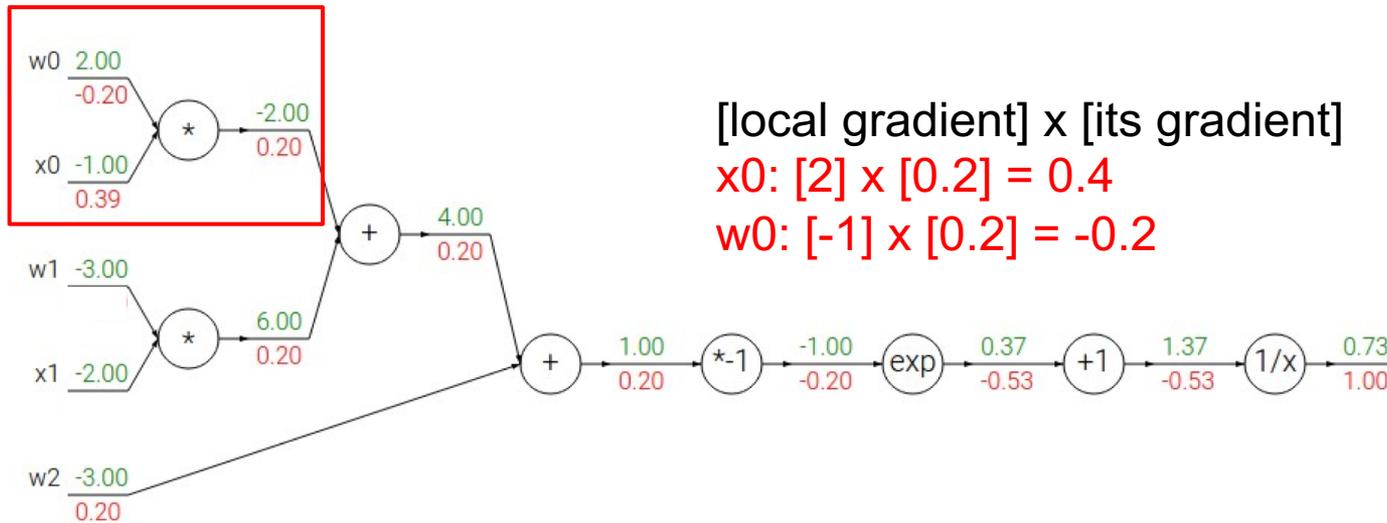
Green- forward computation  
Red – backward derivatives



$f(x) = e^x$	→	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	→	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	→	$\frac{df}{dx} = a$		$f_c(x) = c + x$	→	$\frac{df}{dx} = 1$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$$



[local gradient] x [its gradient]

$x_0: [2] \times [0.2] = 0.4$

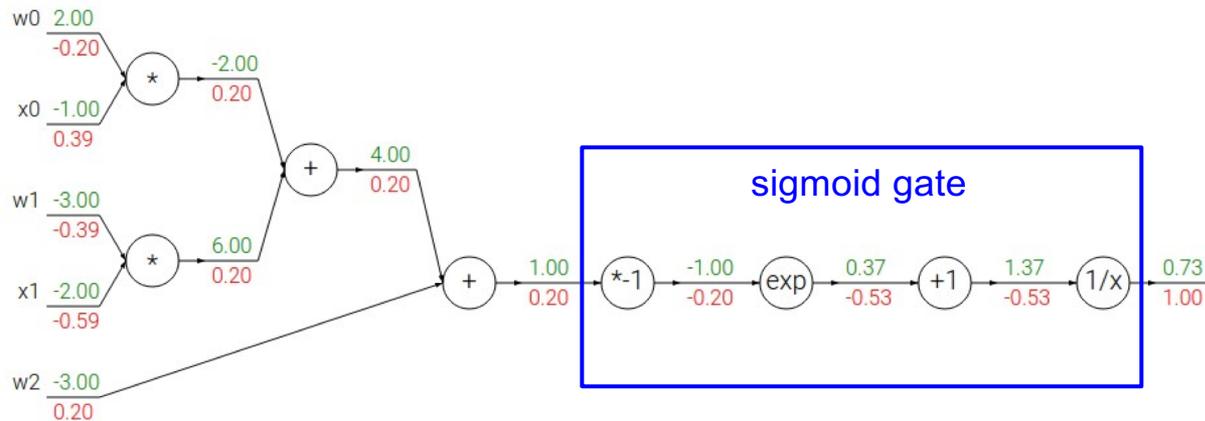
$w_0: [-1] \times [0.2] = -0.2$

$f(x) = e^x$	$\rightarrow$	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	$\rightarrow$	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	$\rightarrow$	$\frac{df}{dx} = a$		$f_c(x) = c + x$	$\rightarrow$	$\frac{df}{dx} = 1$

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \text{sigmoid function}$$

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$



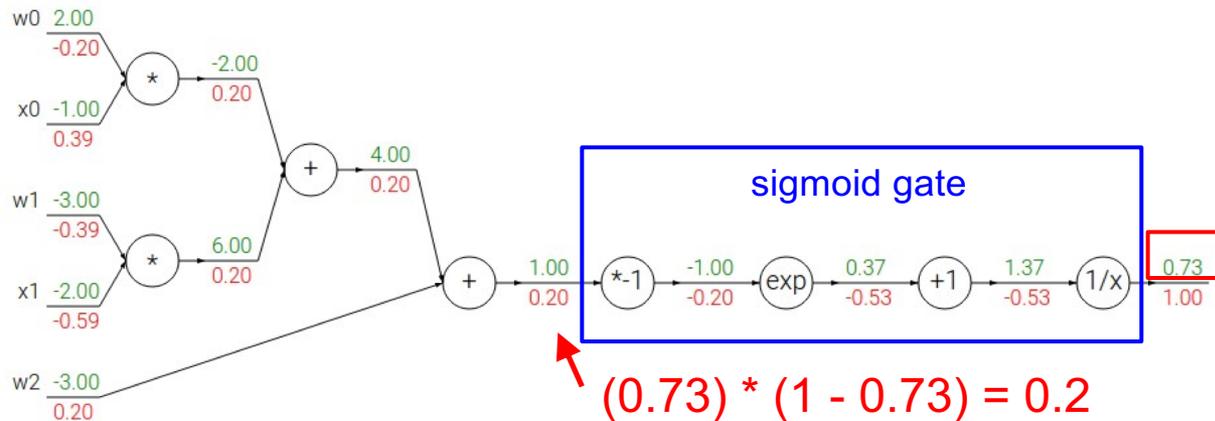
56

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$$

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

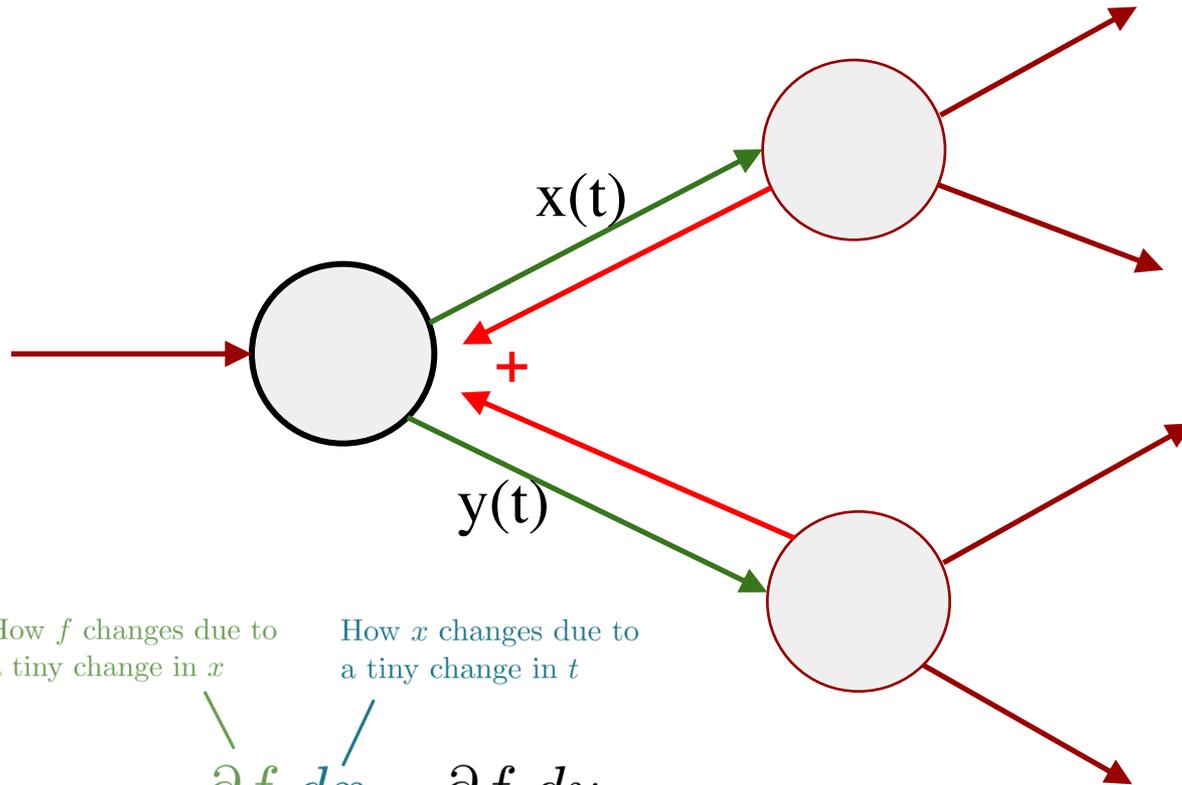
sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$



Green- forward computation  
Red – backward derivatives

# Gradients add at branches



How  $f$  changes due to a tiny change in  $x$

How  $x$  changes due to a tiny change in  $t$

$$\frac{d}{dt} f(x(t), y(t)) = \underbrace{\frac{\partial f}{\partial x} \frac{dx}{dt}}_{\text{Total change in } f \text{ due to the influence } t \text{ has on } x} + \underbrace{\frac{\partial f}{\partial y} \frac{dy}{dt}}_{\text{Total change in } f \text{ due to the influence } t \text{ has on } y}$$

58

This is an ordinary derivative not a partial derivative  $\frac{\partial}{\partial t}$ , because the total composition has one input and one output.

# Summary

---

- SGD
  - Simple linear classifier
  - Complex classification prediction functions
- Computing partial derivatives algorithmically
  - Forward propagation to compute intermediate function values
  - Backward propagation to compute derivatives
- Deep learning
  - New direction for text data processing given its success in image/audio processing
  - Frameworks and software
    - TensorFlow (Google).
    - Others: Theano, Torch, CAFFE, computation graph toolkit (CGT)