# Машинное обучение

на примере глубокого обучения
в компьютерного зрения

## Занятие 3
## Backpropagation и Нейронные сети

Дмитрий Яшунин, к.ф.-м.н
IntelliVision

e-mail: yashuninda@yandex.ru

# На прошлом занятии: **Классификация изображений**

Ключевая задача компьютерного зрения



К какому классу принадлежит изображение?
классы: человек, животное, автомобиль …

кот

# На прошлом занятии: Линейный классификатор

Image



Array of **32x32x3** numbers
(3072 numbers total)

**CIFAR-10**
**50,000** training images
**10,000** testing images
**10** classes

$$s=\underset{10\times1}{f(x,W)}=\underset{10\times3072}{W}\underset{3072\times1}{x}+\underset{10\times1}{b}$$

s – scores
W – weights or parameters
x – image pixels
b – bias

# На прошлом занятии: Интерпретация линейного классификатора

**CIFAR-10**



$$f(x,W) = Wx + b$$



Example trained weights of a linear classifier trained on CIFAR-10:

# На прошлом занятии: Функции ошибки

Image



$x_i$ - image

$y_i$ - label, element of a set {0, 1, ...}

scores $s = f(x_i, W) = [s_0, \dots s_{y_i}, \dots]$

**Loss over dataset:**

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i$$

**Multiclass SVM (hinge) loss:**

$$L_i = \sum_{i \neq y_i} \max(0, s_i - s_{y_i} + 1)$$

**Cross-entropy (softmax) loss:**

$$L_i = -\log \frac{e^{s_{y_i}}}{\sum_j e^{s_j}}$$

# На прошлом занятии: Регуляризация

Softmax or SVM

λ - regularization strength (hyperparameter)

Full loss

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda R(W)$$

Data loss    Regularization

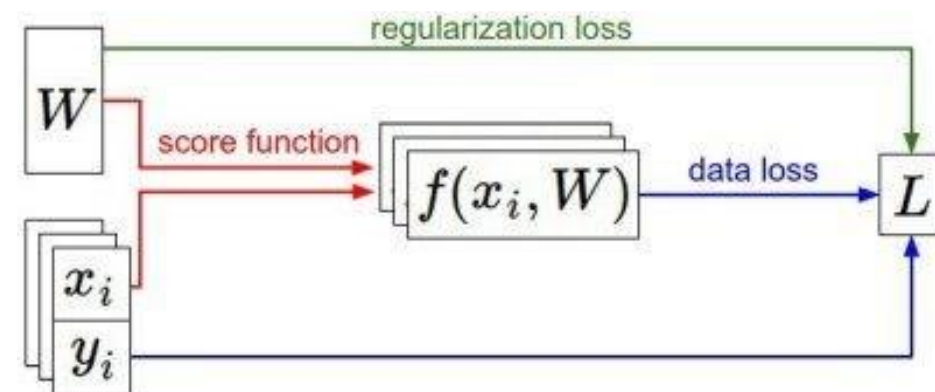How do we find the best W?

**L2 regularization** $\quad R(W) = \sum_k \sum_l W_{k,l}^2$
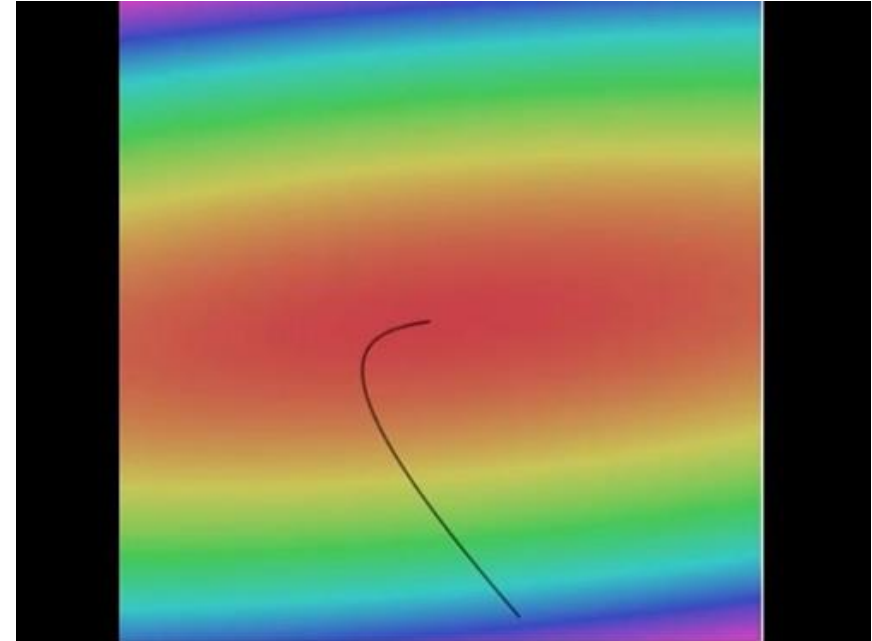
L1 regularization $\quad R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2) $\quad R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

# Оптимизация





```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```
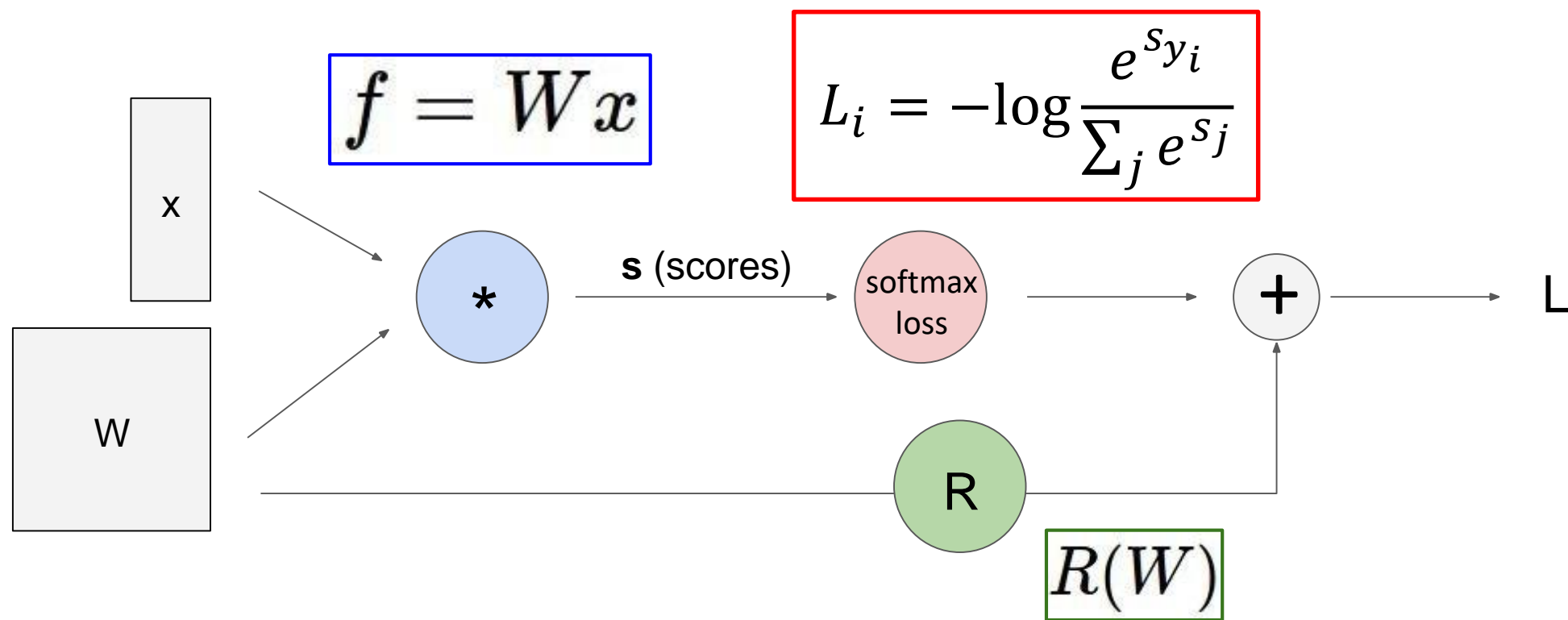
# Метод градиентного спуска

$$\frac{dL(w)}{dw} = \lim_{h \to 0} \frac{L(w+h) - L(w)}{h}$$

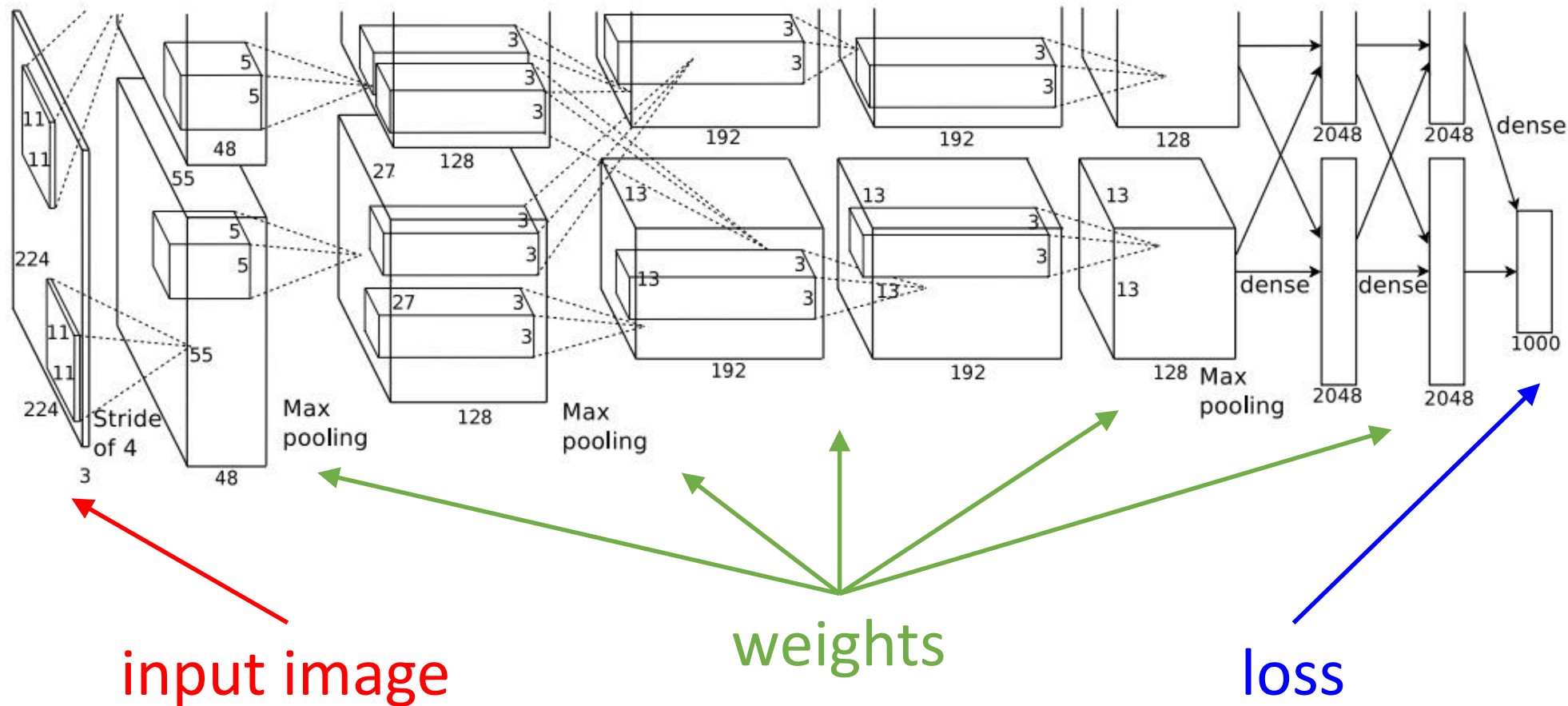**Численные градиенты**: медленно, не точно, быстро реализовать

**Аналитические градиенты**: быстро, точно, можно ошибиться

# Вычислительный граф



$$f = Wx$$

$$L_i = -\log \frac{e^{s_{y_i}}}{\sum_j e^{s_j}}$$

x

W

\* 

**s** (scores)

softmax loss

R

+

L

$$R(W)$$

# Вычислительный граф: Convolutional Network

## AlexNet



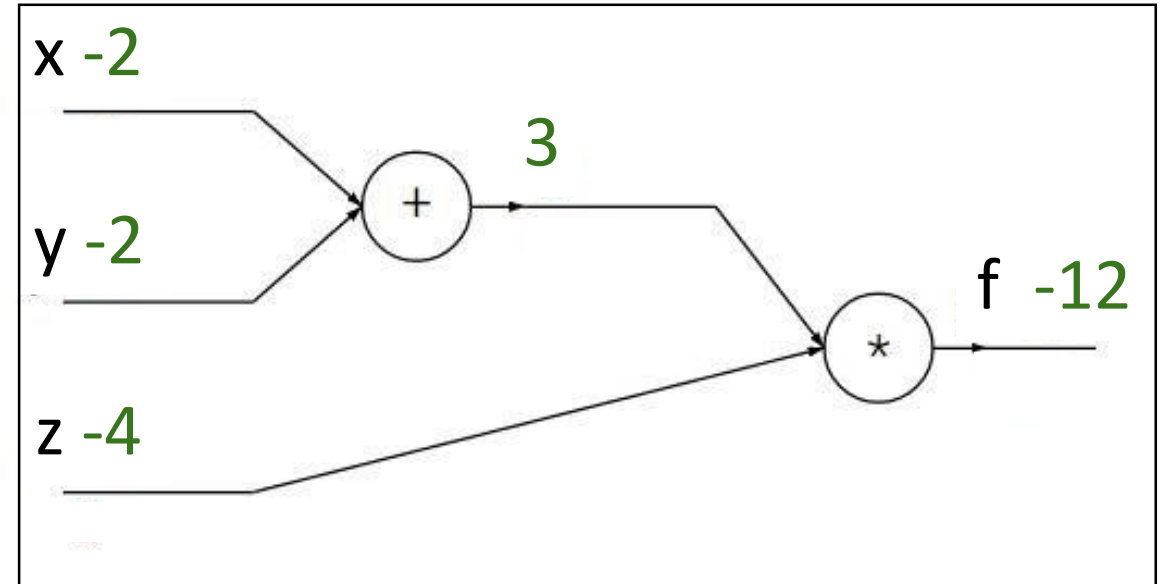input image

weights

loss

# Backpropagation

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
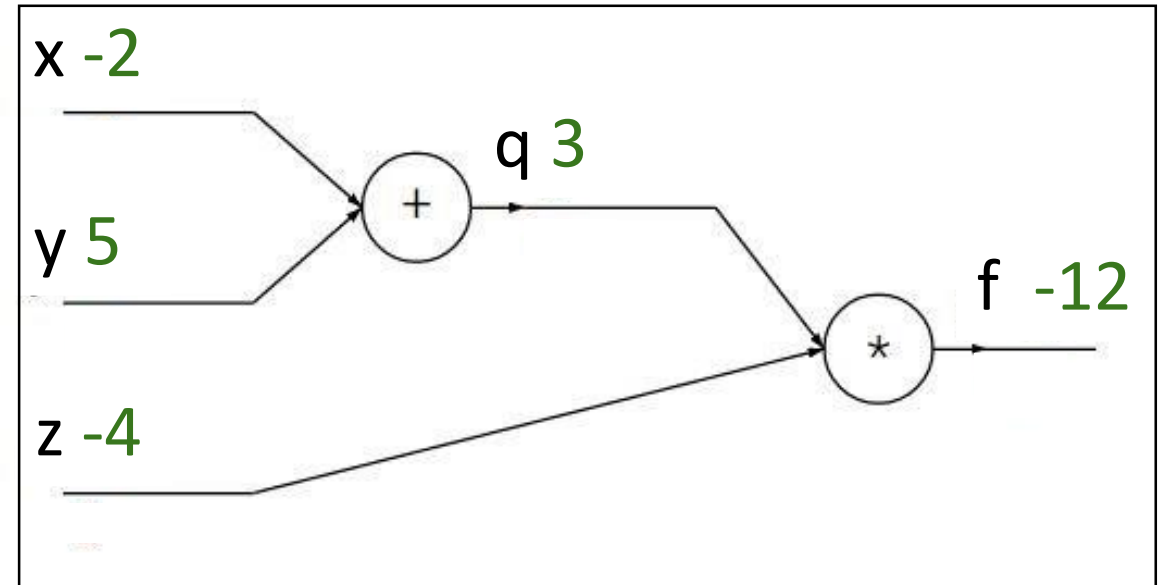
# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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

x -2

q 3

y 5

f -12

z -4

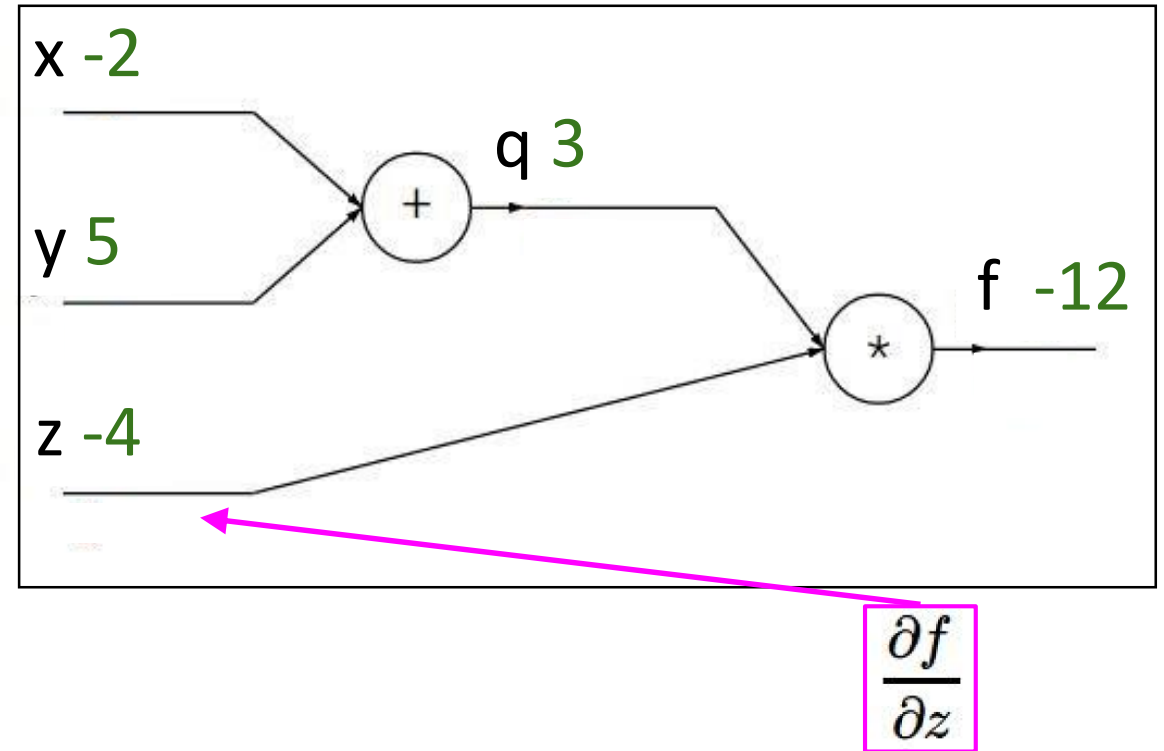$$\frac{\partial f}{\partial z}$$

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
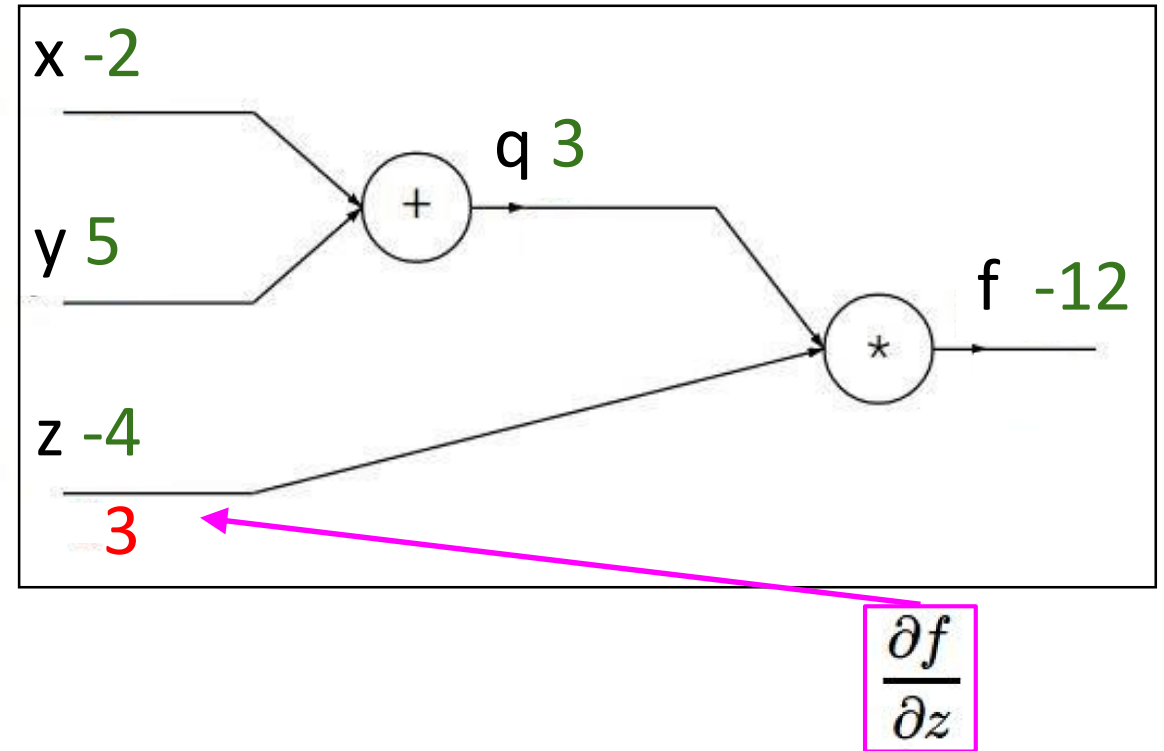
# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
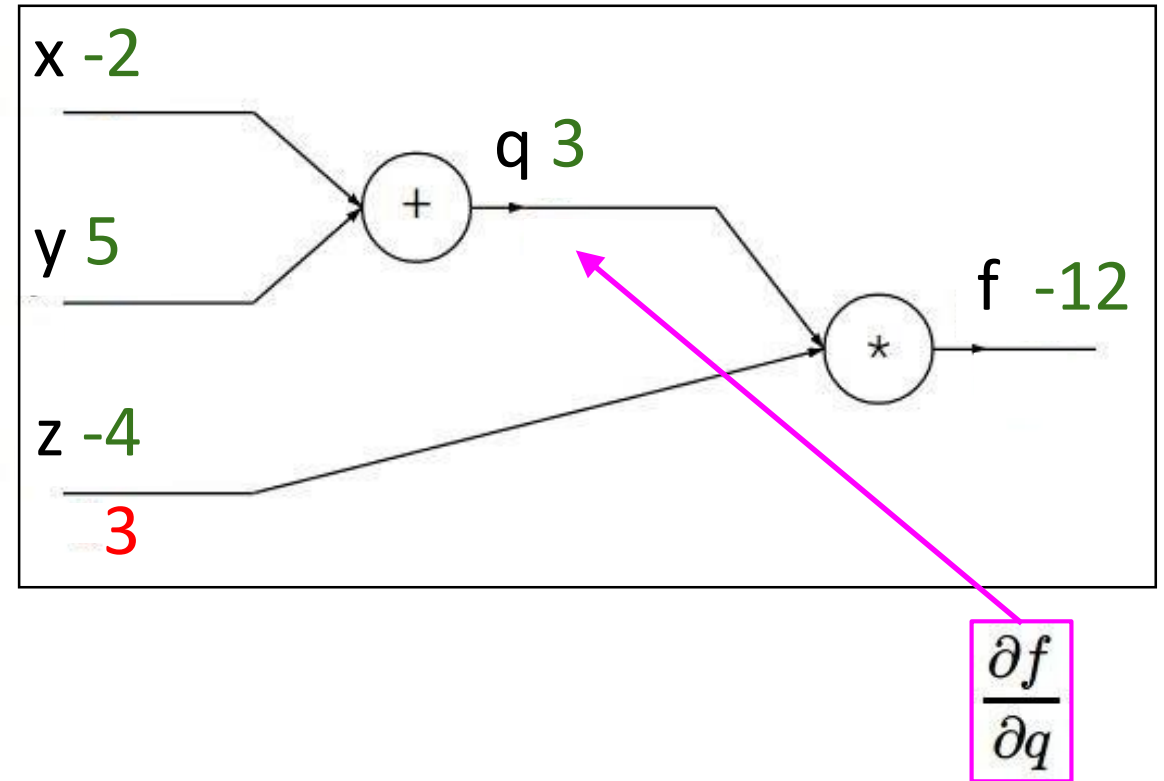
# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
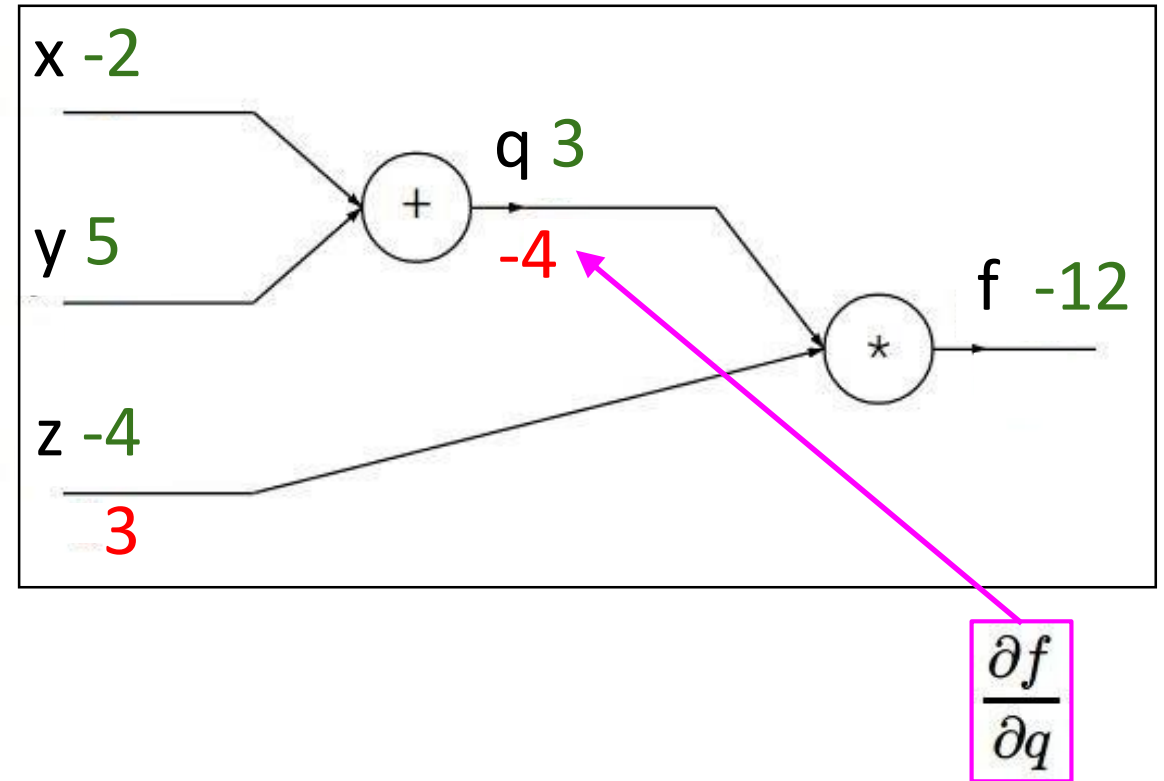
# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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



x -2

q 3

y 5

-4

f -12

z -4

3

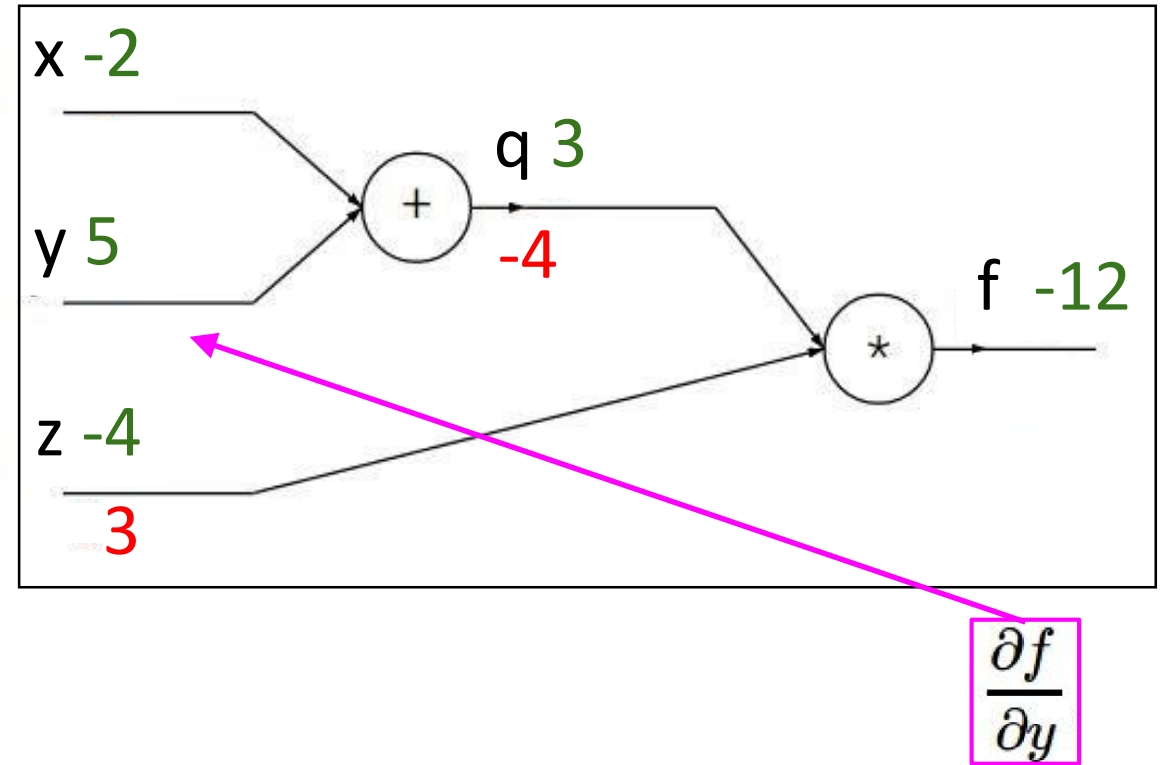$\dfrac{\partial f}{\partial y}$

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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



x -2

q 3

y 5

-4

f -12

z -4

3

Правило диффиренцирования:

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

$\dfrac{\partial f}{\partial y}$

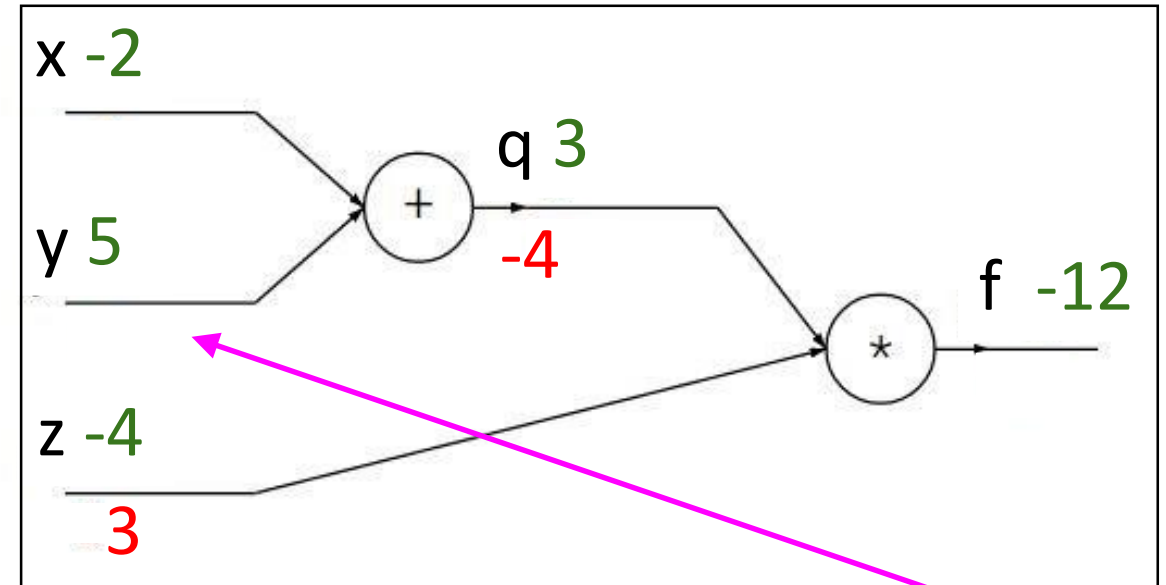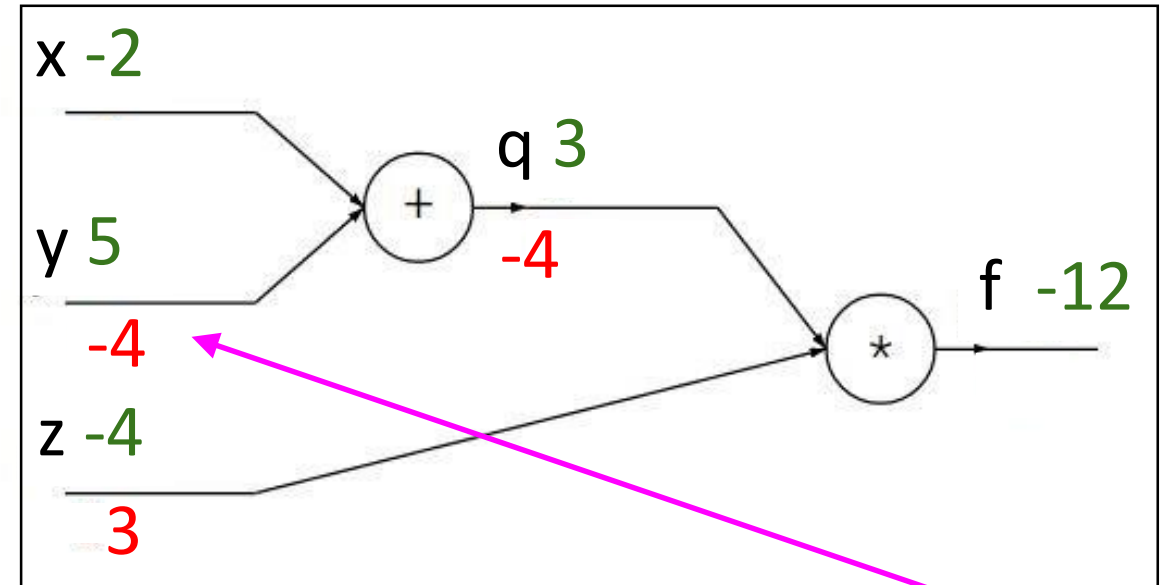Входящий градиент

Локальный градиент

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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

x -2

y 5

q 3

-4

z -4

-4

3

f -12

$\dfrac{\partial f}{\partial y}$

Правило диффиренцирования:

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

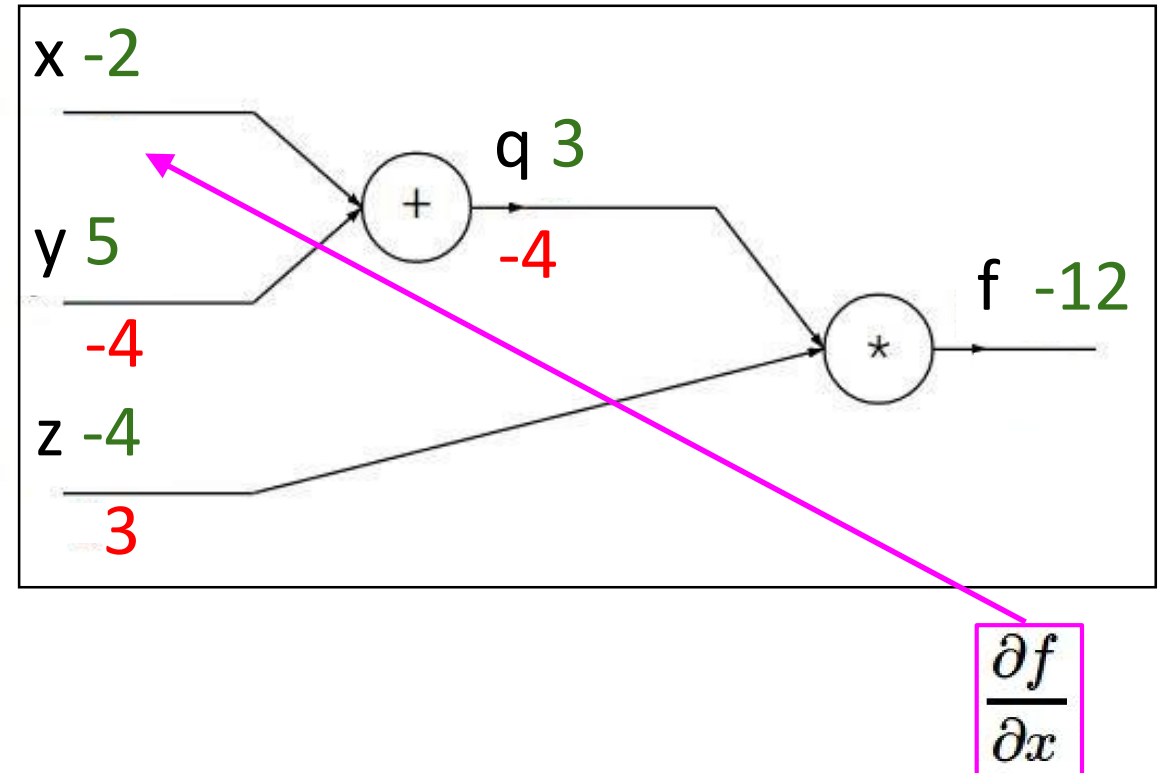Входящий градиент

Локальный градиент

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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



x -2

q 3

y 5

-4

-4

z -4

3

f -12

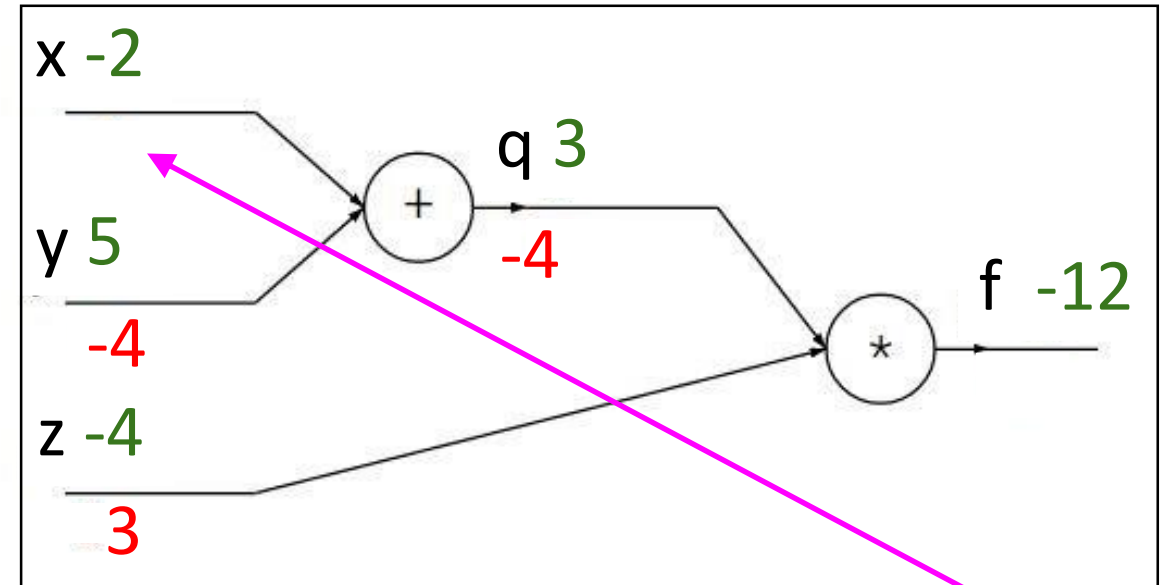$$\frac{\partial f}{\partial x}$$

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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

x -2

y 5
-4

z -4
3

q 3
-4

f -12

Правило диффиренцирования:

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

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

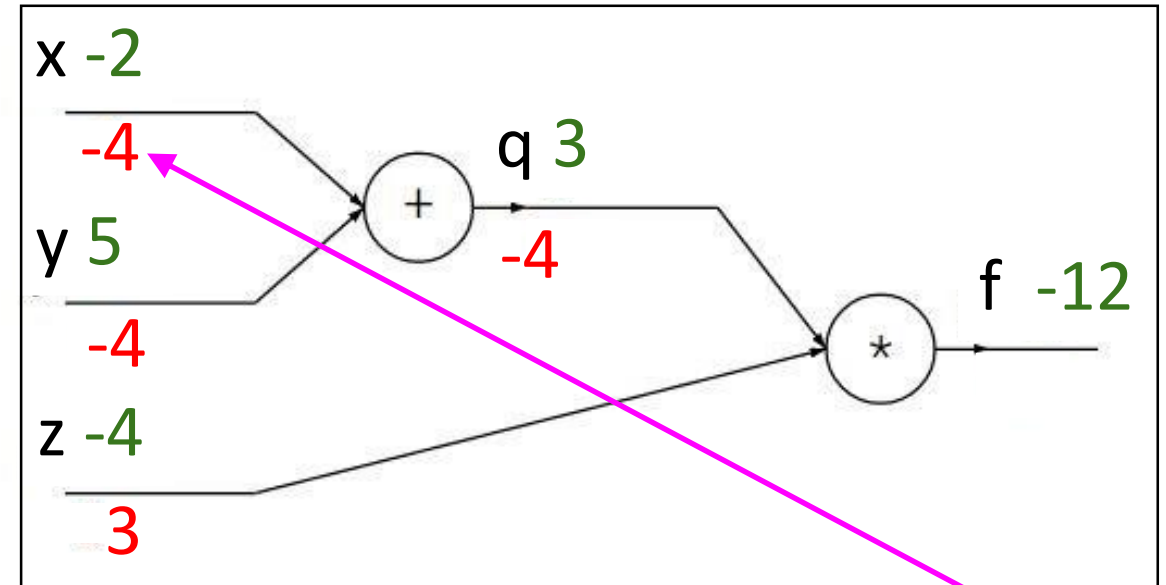Входящий градиент

Локальный градиент

# Backpropagation: пример

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

e.g. x = -2, y = 5, z = -4

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

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

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

x -2

-4

q 3

y 5

-4

-4

z -4

3

f -12

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

Правило диффиренцирования:
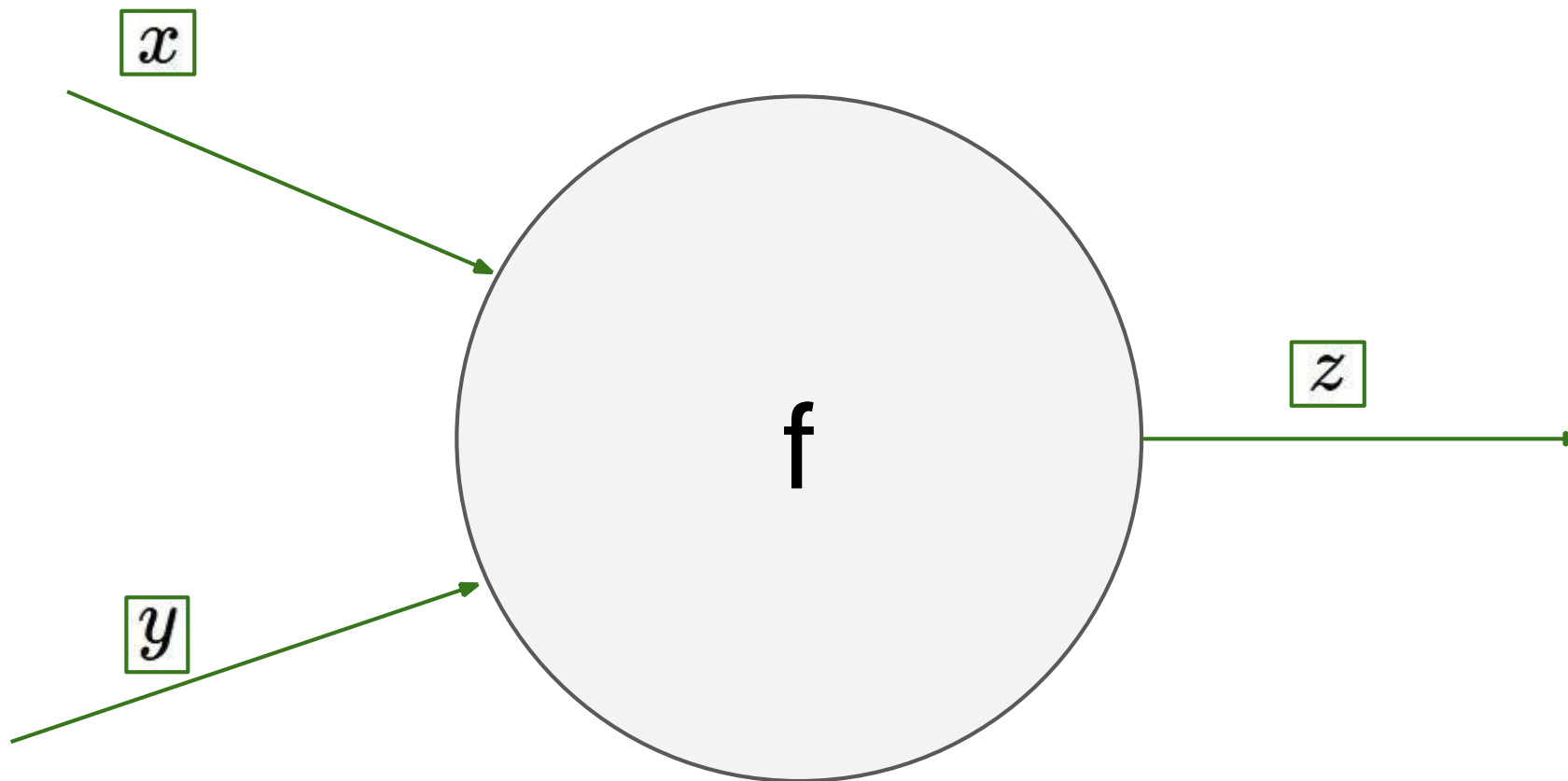
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$
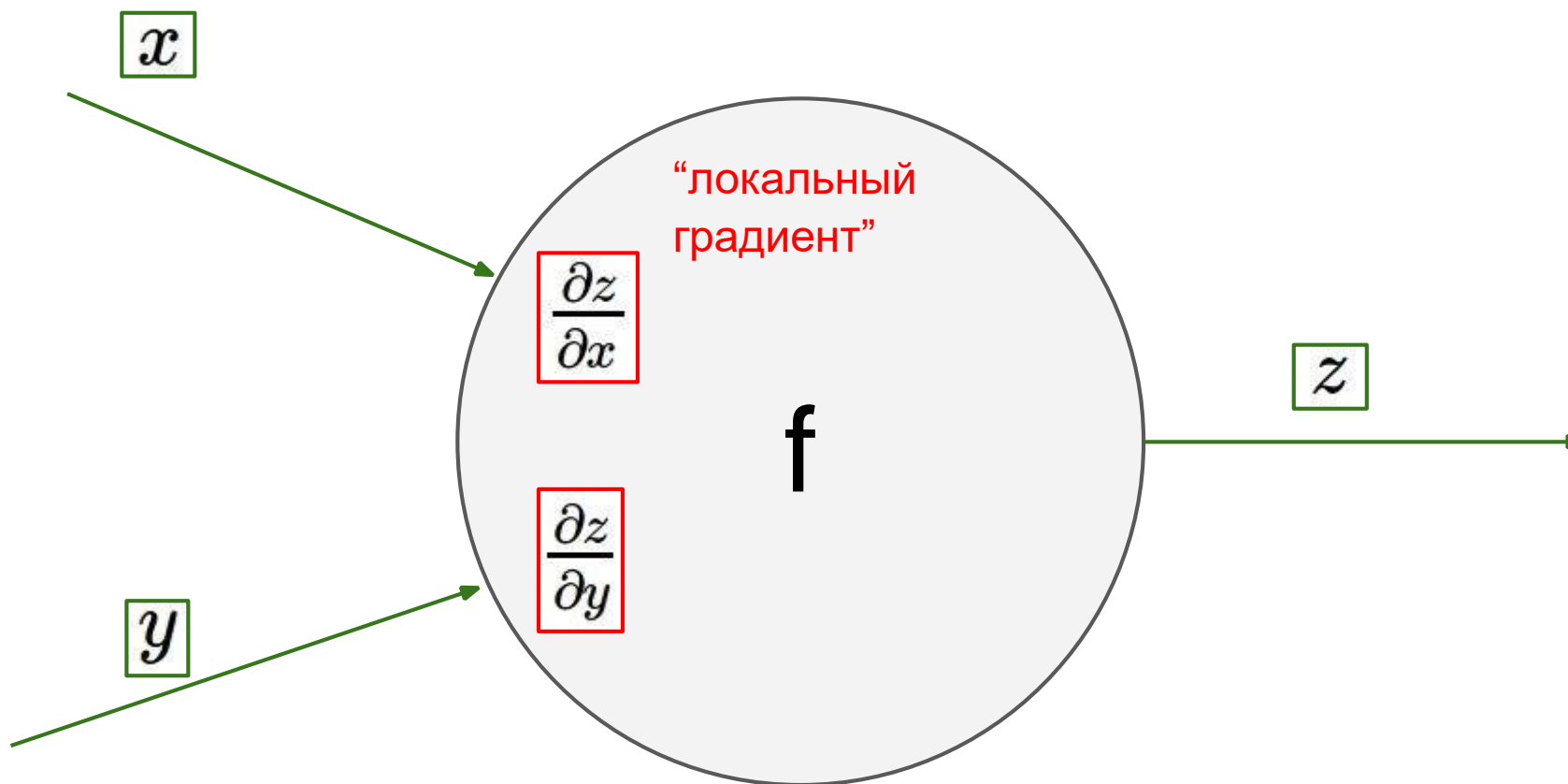
Входящий градиент
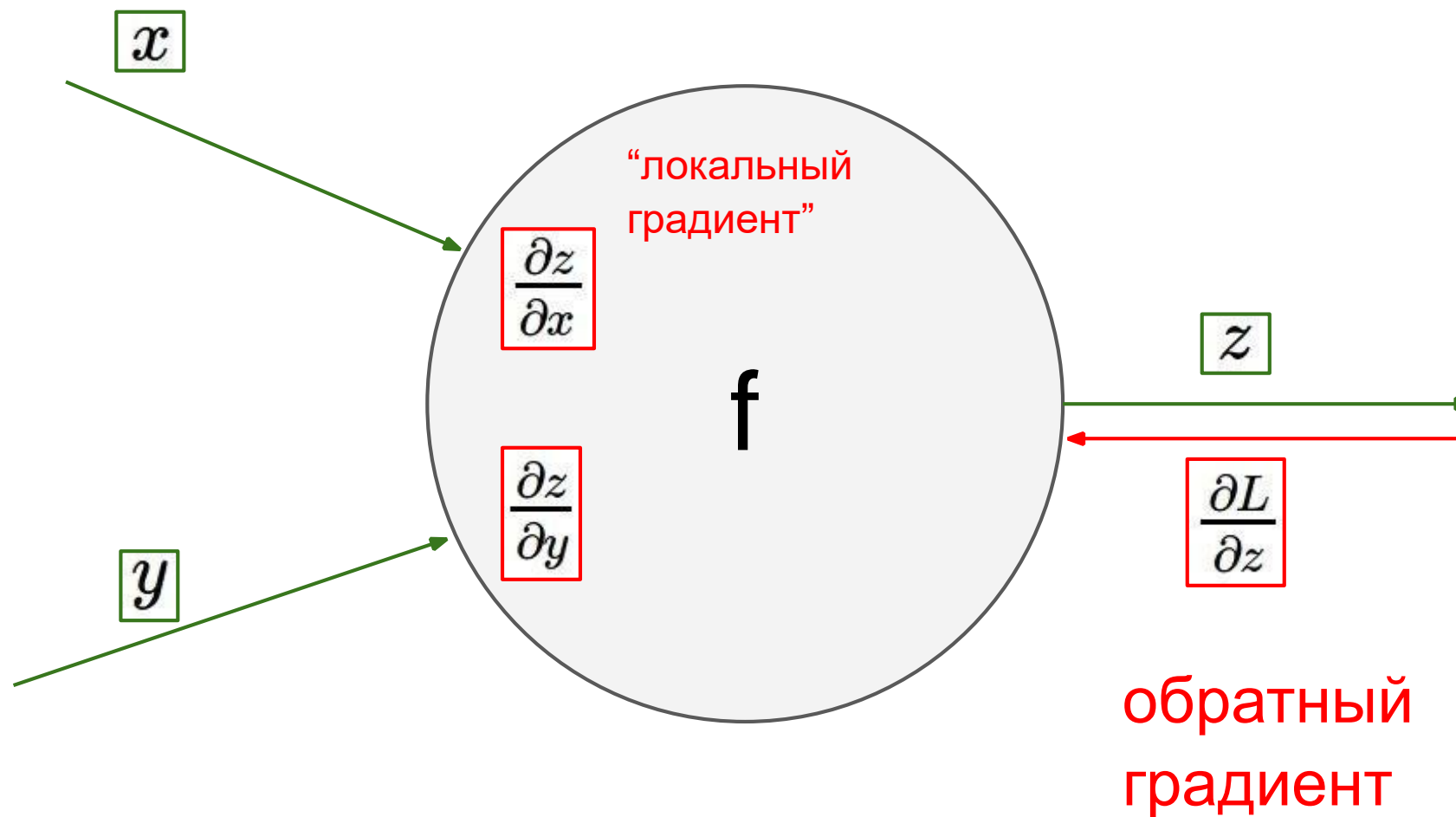
Локальный градиент

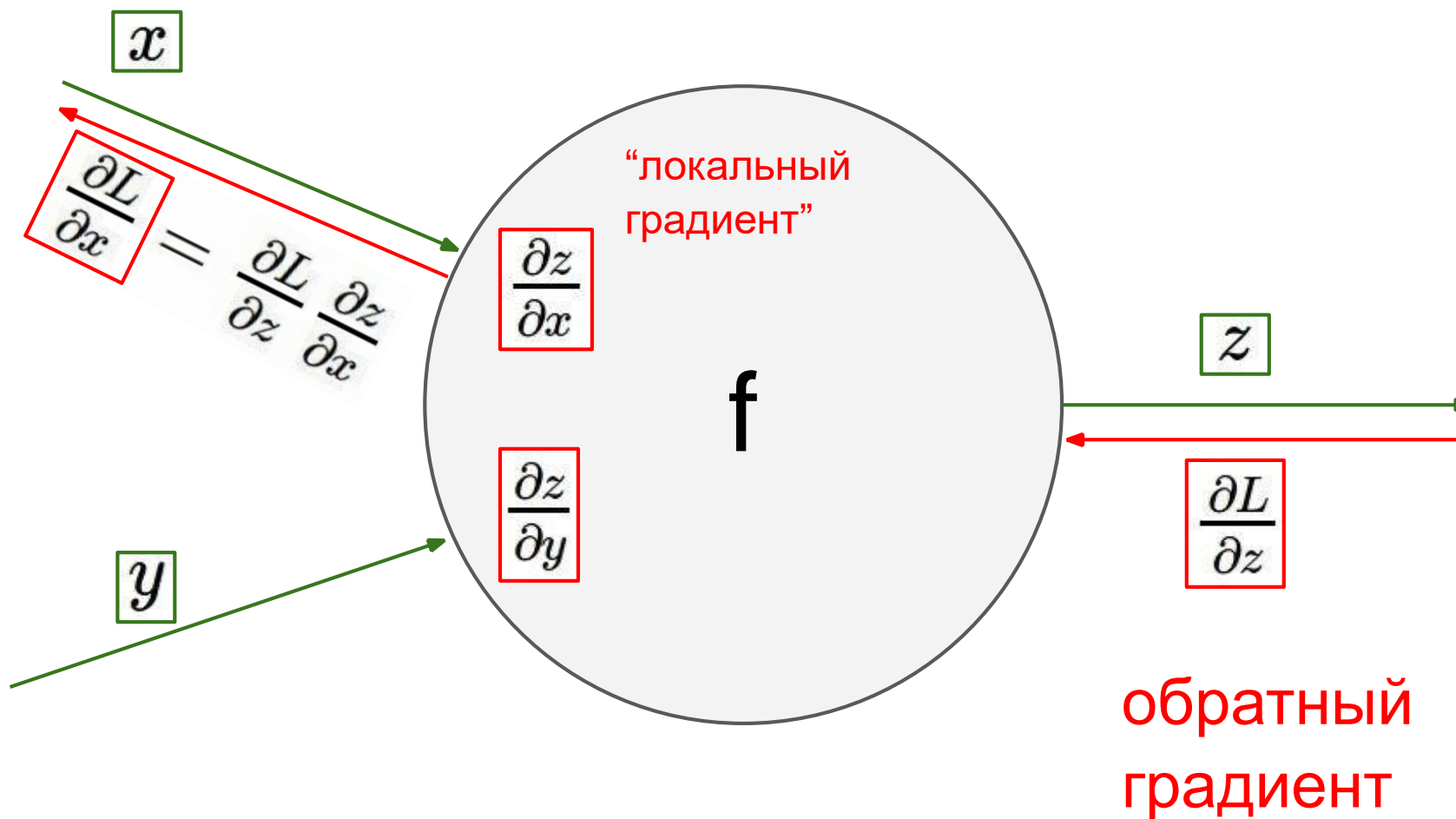# Backpropagation: общий случай

# Backpropagation: общий случай

# Backpropagation: общий случай

# Backpropagation: общий случай
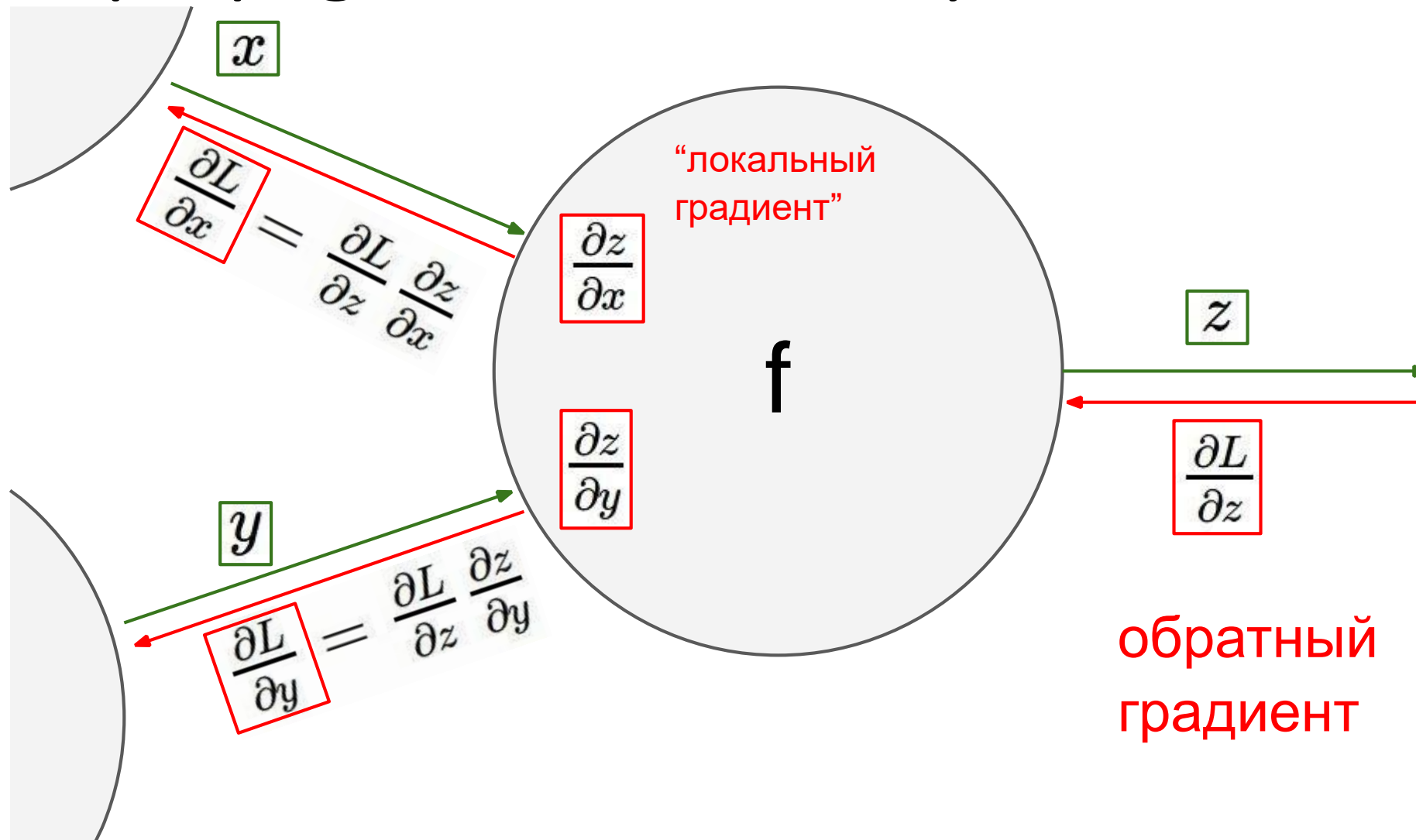


$x$

$$\boxed{\frac{\partial L}{\partial x}} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

"локальный градиент"

$\boxed{\dfrac{\partial z}{\partial x}}$

$\boxed{\dfrac{\partial z}{\partial y}}$

f

$y$

$z$

$\boxed{\dfrac{\partial L}{\partial z}}$

обратный градиент

# Backpropagation: общий случай



$x$

$\dfrac{\partial L}{\partial x} = \dfrac{\partial L}{\partial z}\dfrac{\partial z}{\partial x}$

"локальный градиент"

$\dfrac{\partial z}{\partial x}$

**f**

$\dfrac{\partial z}{\partial y}$

$z$

$\dfrac{\partial L}{\partial z}$

$y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial L}{\partial z}\dfrac{\partial z}{\partial y}$

обратный градиент

# Backpropagation: общий случай

# Пример:

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



w0 2.00

-2.00

x0 -1.00

w1 -3.00

6.00

4.00

x1 -2.00

w2 -3.00

1.00   -1.00   0.37   1.37   0.73

# Пример:

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

w0 2.00

-2.00

x0 -1.00

4.00

w1 -3.00

6.00

1.00    -1.00    0.37    1.37    0.73

x1 -2.00    *    +    *-1    exp    +1    1/x

w2 -3.00

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

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Пример:

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

w0 2.00

x0 -1.00

-2.00

w1 -3.00

4.00

6.00

x1 -2.00

1.00    -1.00    0.37    1.37    0.73

1.00

w2 -3.00

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

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Пример:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$
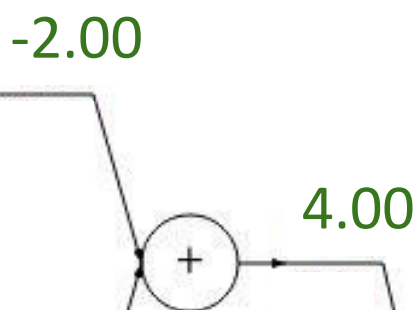
w0 2.00

-2.00

x0 -1.00

4.00

w1 -3.00

6.00

x1 -2.00

1.00   -1.00   0.37   1.37   0.73

w2 -3.00

Локальный градиент

Входящий градиент

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

-0.53   1.00

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

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Пример:

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

w0 2.00

x0 -1.00

w1 -3.00

x1 -2.00

w2 -3.00



-2.00

4.00

6.00

1.00    -1.00    0.37    1.37    0.73

-0.53    1.00

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

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

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

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

# Пример:

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

w0 2.00

x0 -1.00

-2.00

w1 -3.00

6.00

x1 -2.00

4.00

w2 -3.00

1.00   -1.00

0.37   1.37   0.73

-0.53   -0.53   1.00

Локальный градиент

Входящий градиент

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

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

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

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

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

# Пример:

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



w0 2.00

-2.00

x0 -1.00

w1 -3.00

6.00

4.00

x1 -2.00

w2 -3.00

1.00    -1.00    0.37    1.37    0.73

exp    +1    1/x

-0.53    -0.53    1.00

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

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

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

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

# Пример:

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

w0 2.00

-2.00

x0 -1.00

w1 -3.00

4.00

6.00

x1 -2.00

w2 -3.00

Локальный градиент

Входящий градиент

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

1.00

-1.00          0.37

1.37          0.73

-0.20          -0.53

-0.53          1.00

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

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

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

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

# Пример:

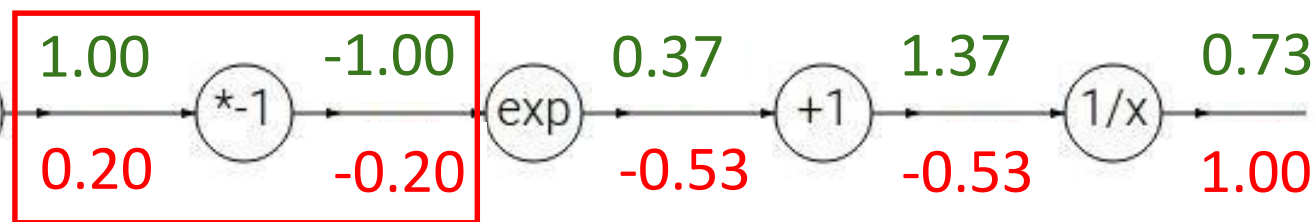$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



w0 2.00

x0 -1.00          -2.00

w1 -3.00

x1 -2.00          6.00          4.00

                  1.00          -1.00          0.37          1.37          0.73

                                -0.20          -0.53          -0.53          1.00

w2 -3.00

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

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

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

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

# Пример:

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

w0 2.00

x0 -1.00

-2.00

w1 -3.00

6.00

x1 -2.00

4.00

w2 -3.00

Локальный градиент    Входящий градиент

(-1)*0.20=0.20

1.00    -1.00

0.20    -0.20

0.37

-0.53

1.37

-0.53

0.73

1.00

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

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

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

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

# Пример:

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

w0 2.00

x0 -1.00

-2.00

w1 -3.00

6.00

x1 -2.00

w2 -3.00

4.00

1.00
0.20

-1.00
-0.20

0.37
-0.53

1.37
-0.53

0.73
1.00

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

$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$

# Пример:

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



w0 2.00

x0 -1.00

-2.00

w1 -3.00

x1 -2.00

6.00

w2 -3.00

0.2

4.00

0.2

1.00

0.20

*-1

-1.00

-0.20

exp

0.37

-0.53

+1

1.37

-0.53

1/x

0.73

1.00

[локальный градиент] x [входящий градиент]

[1] x [0.2] = 0.2

[1] x [0.2] = 0.2

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

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

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

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

# Пример:

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

w0 2.00

-2.00

x0 -1.00

0.2

w1 -3.00

6.00

x1 -2.00

0.2

4.00

0.2

w2 -3.00

0.2

1.00

0.20

-1.00

-0.20

0.37

-0.53

1.37

-0.53

0.73

1.00

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

$f_a(x) = ax$ $\rightarrow$ $\frac{df}{dx} = a$ $\quad\bigg|\quad$ $f_c(x) = c + x$ $\rightarrow$ $\frac{df}{dx} = 1$

# Пример:

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



[локальный градиент] x [входящий градиент]

x0: [2] x [0.2] = 0.4

w0: [-1] x [0.2] = -0.2

w0 2.00
-0.2
x0 -1.00
0.4

-2.00
0.2

w1 -3.00
6.00
x1 -2.00
0.2

4.00
0.2

1.00        -1.00        0.37        1.37        0.73
0.20        -0.20        -0.53       -0.53       1.00

w2 -3.00
0.2

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \quad \bigg| \quad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \quad \bigg| \quad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

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

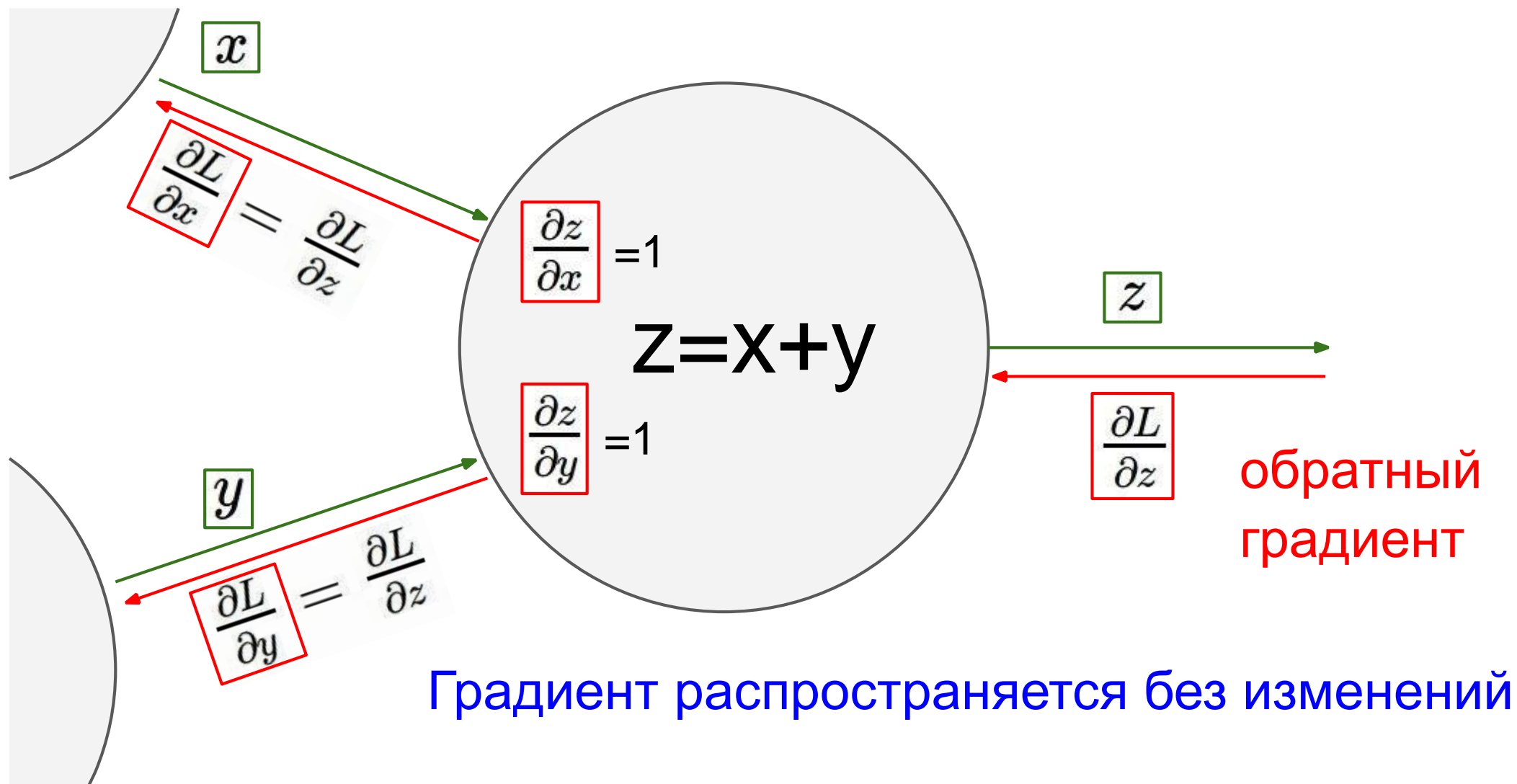$$\boxed{\sigma(x) = \frac{1}{1 + e^{-x}}}$$ sigmoid function

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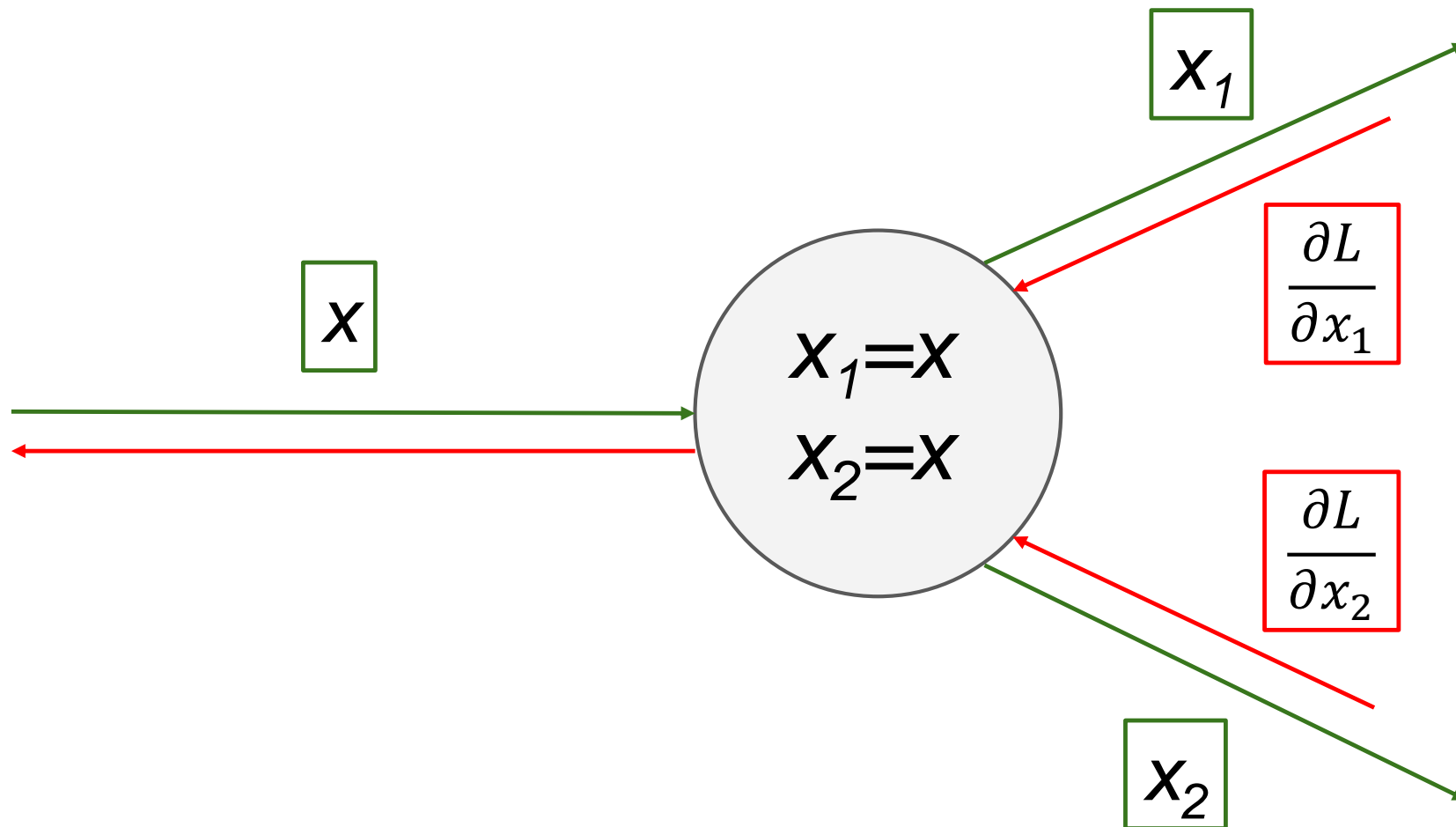


w0 2.00
-0.2
x0 -1.00
0.4
-2.00
0.2

w1 -3.00
6.00
0.2
x1 -2.00
0.2

4.00
0.2

w2 -3.00
0.2

1.00
0.20

sigmoid function

-1.00
-0.20

0.37
-0.53

1.37
-0.53

0.73
1.00

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

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

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



w0 2.00
-0.2
x0 -1.00
0.4
-2.00
0.2

w1 -3.00
6.00
x1 -2.00
0.2
4.00
0.2

w2 -3.00
0.2

1.00
0.20

sigmoid function

-1.00
-0.20

0.37
-0.53

1.37
-0.53

0.73
1.00

$*{-1}$  exp  $+1$  $1/x$

(0.73) * (1 - 0.73) = 0.2

# Обратный градиент при суммировании

# Обратный градиент при суммировании



$$\boxed{\frac{\partial L}{\partial x}} = \frac{\partial L}{\partial z}$$

$x$

$\boxed{\frac{\partial z}{\partial x}} = 1$

$z = x + y$

$\boxed{\frac{\partial z}{\partial y}} = 1$

$z$

$$\boxed{\frac{\partial L}{\partial y}} = \frac{\partial L}{\partial z}$$

$y$

$\boxed{\frac{\partial L}{\partial z}}$

обратный градиент

Градиент распространяется без изменений

# Обратный градиент при переиспользовании переменной

# Обратный градиент при переиспользовании переменной



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial x_1}\frac{\partial x_1}{\partial x} + \frac{\partial L}{\partial x_2}\frac{\partial x_2}{\partial x_1}$$

$x$

$x_1=x$
$x_2=x$

$x_1$

$\frac{\partial L}{\partial x_1}$

$\frac{\partial L}{\partial x_2}$

$x_2$

# Обратный градиент при переиспользовании переменной



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial x_1} + \frac{\partial L}{\partial x_2}$$

Градиенты суммируются

$x_1=x$
$x_2=x$

$x$

$x_1$

$x_2$

$\frac{\partial L}{\partial x_1}$

$\frac{\partial L}{\partial x_2}$

# Backpropagation: векторный случай

(x,y,z are now vectors)

This is now the **Jacobian matrix** (derivative of each element of z w.r.t. each element of x)

"local gradient"

$$\boxed{\frac{\partial L}{\partial x}} = \frac{\partial z}{\partial x} \frac{\partial L}{\partial z}$$

$\boxed{x}$

$\boxed{\frac{\partial z}{\partial x}}$

$\boxed{\frac{\partial z}{\partial y}}$

$\boxed{y}$

f

$\boxed{z}$

$\boxed{\frac{\partial L}{\partial z}}$

gradients

# Векторные операции

4096-d
input vector

$$f(x) = \max(0, x)$$
*(elementwise)*

4096-d
output vector

# Векторные операции

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

Q: what is the
size of the
Jacobian matrix?

# Векторные операции

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

# Векторные операции

4096-d
input vector

$f(x) = \max(0,x)$
*(elementwise)*

4096-d
output vector

Q: what is the size of the Jacobian matrix? [4096 x 4096!]

in practice we process an entire minibatch (e.g. 100) of examples at one time:

i.e. Jacobian would technically be a [409,600 x 409,600] matrix

# Векторные операции

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Q2: what does it
look like?

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$$

$$\in \mathbb{R}^n \in \mathbb{R}^{n \times n}$$

W

x

*

L2

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$$

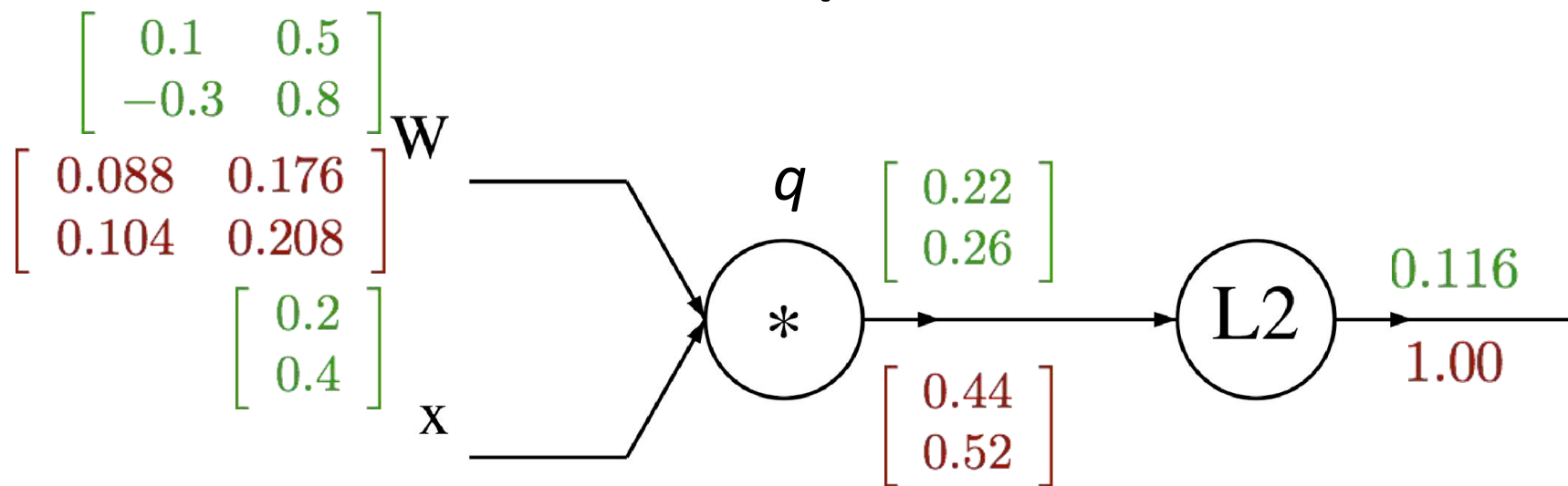$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$$

$q$

$*$ → L2 →

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$$

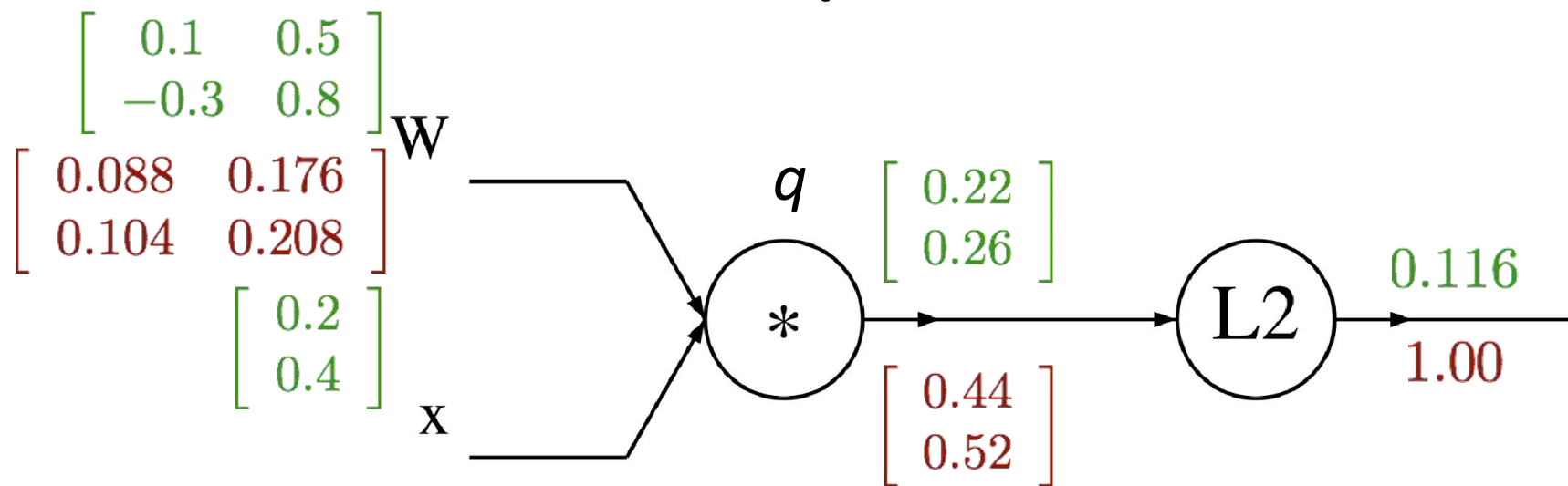$q$ $\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

\*

0.116

L2

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$

$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$

$q$ $\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

$*$

L2

0.116

1.00

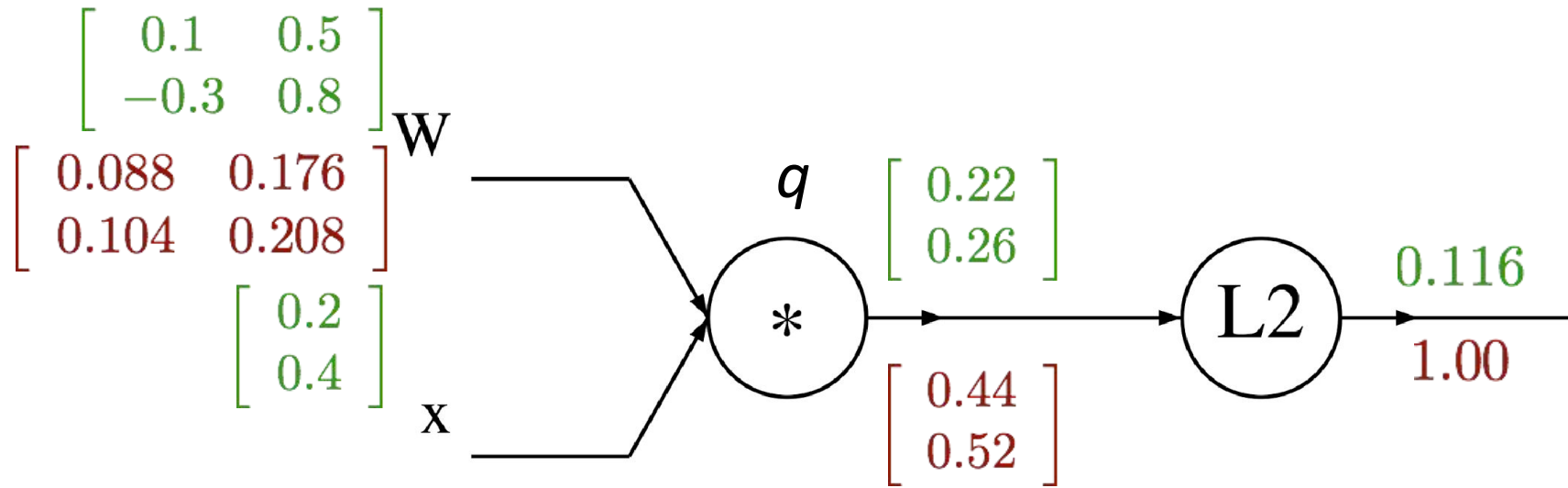$q = W \cdot x = \begin{pmatrix} W_{1,1} x_1 + \cdots + W_{1,n} x_n \\ \vdots \\ W_{n,1} x_1 + \cdots + W_{n,n} x_n \end{pmatrix}$

$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$$

$q$   $\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

$*$     L2   0.116 / 1.00

$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\nabla_q f = 2q$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$$

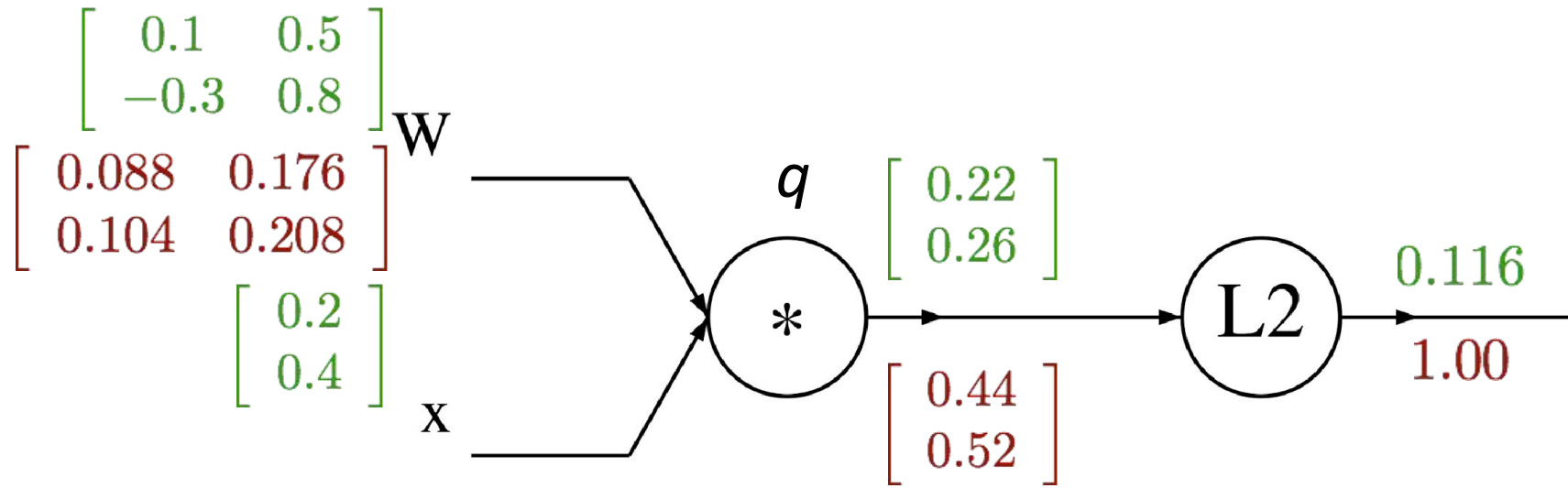$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$$

$q$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$*$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\nabla_q f = 2q$$
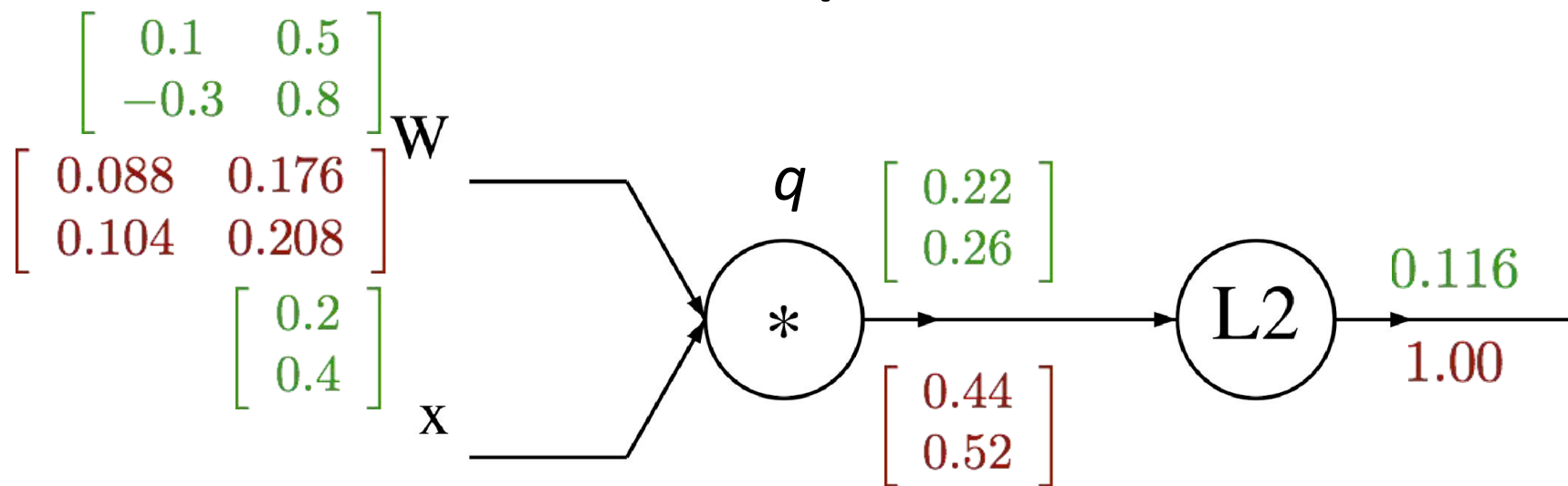
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:    $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$

$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$

$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$

$q$

$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

$*$

$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$

L2

0.116

1.00

$\dfrac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$

$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$
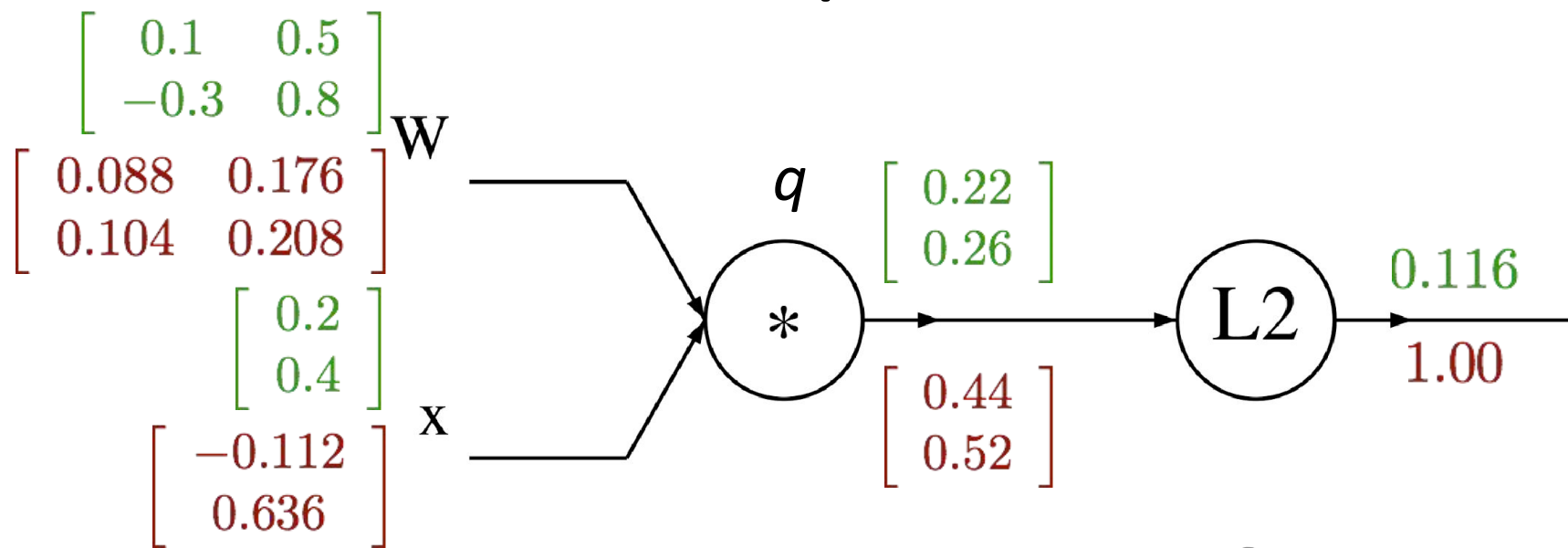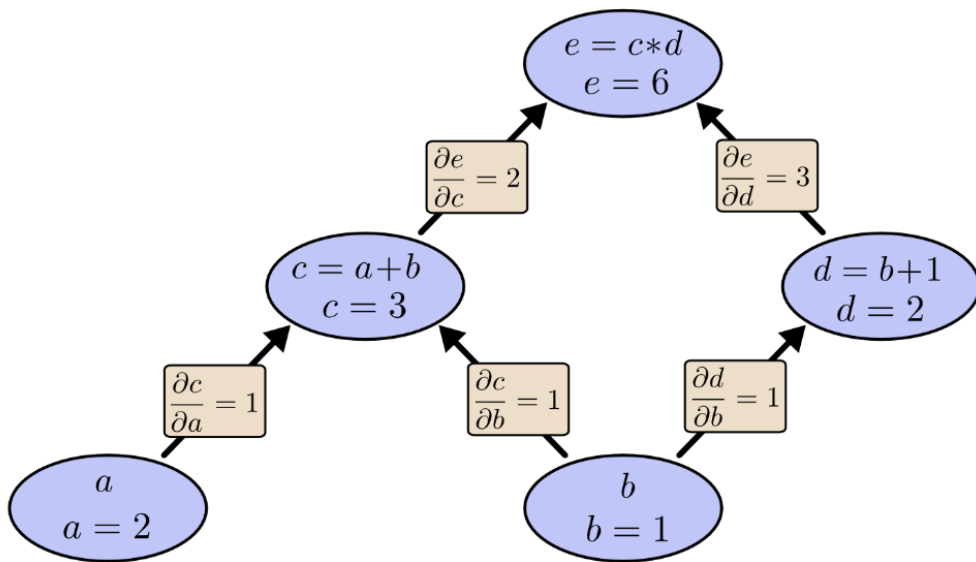
$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_W$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_x$$

$q$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$*$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

$q$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

*

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1} x_1 + \cdots + W_{1,n} x_n \\ \vdots \\ W_{n,1} x_1 + \cdots + W_{n,n} x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\nabla_W f = 2q \cdot x^T$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

$q$ $\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

$x$

$*$  $\longrightarrow$  $L2$  $\xrightarrow{\text{0.116}}$

1.00

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1} x_1 + \cdots + W_{1,n} x_n \\ \vdots \\ W_{n,1} x_1 + \cdots + W_{n,n} x_n \end{pmatrix}$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

$$\nabla_W f = 2q \cdot x^T$$

$q$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

*

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

Всегда проверяйте размерность градиента

Размерность градиента должна быть такой-же как размерность переменной, для которой он вычисляется

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

*

$q$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1} x_1 + \cdots + W_{1,n} x_n \\ \vdots \\ W_{n,1} x_1 + \cdots + W_{n,n} x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

# A vectorized example:  $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}$ W

$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$

$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$ x

$q$ $\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

$*$

$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$

L2

0.116

1.00

$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$

$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$

$\dfrac{\partial q_k}{\partial x_i} = W_{k,i}$

$\dfrac{\partial f}{\partial x_i} = \sum_k \dfrac{\partial f}{\partial q_k} \dfrac{\partial q_k}{\partial x_i}$

$\dfrac{\partial f}{\partial x_i} = \sum_k 2q_k W_{k,i}$

# A vectorized example:

$$f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \mathrm{W}$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

$$\begin{bmatrix} -0.112 \\ 0.636 \end{bmatrix} \mathrm{x}$$

$q$ $\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$

\* → L2 → 0.116 / 1.00

$$\nabla_x f = 2W^T \cdot q$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

$$\frac{\partial f}{\partial x_i} = \sum_k \frac{\partial f}{\partial q_k}\frac{\partial q_k}{\partial x_i}$$

$$\frac{\partial f}{\partial x_i} = \sum_k 2q_k W_{k,i}$$

# Modularized implementation: forward / backward API

Graph (or Net) object *(rough psuedo code)*



```python
class ComputationalGraph(object):
    #...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

The computational graph diagram:

$$e = c*d$$
$$e = 6$$

$$\frac{\partial e}{\partial c} = 2 \qquad \frac{\partial e}{\partial d} = 3$$

$$c = a+b$$
$$c = 3$$

$$d = b+1$$
$$d = 2$$

$$\frac{\partial c}{\partial a} = 1 \qquad \frac{\partial c}{\partial b} = 1 \qquad \frac{\partial d}{\partial b} = 1$$

$$a$$
$$a = 2$$

$$b$$
$$b = 1$$

# Modularized implementation: forward / backward API



x

z

*

y

(x,y,z are scalars)

```python
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        return z
    def backward(dz):
        # dx = ... #todo
        # dy = ... #todo
        return [dx, dy]
```

$$\frac{\partial L}{\partial z}$$

$$\frac{\partial L}{\partial x}$$

# Modularized implementation: forward / backward API

X

z

*

y

(x,y,z are scalars)

```python
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z
    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

# Example: Caffe layers

# Caffe Sigmoid Layer

```
1    #include <cmath>
2    #include <vector>
3
4    #include "caffe/layers/sigmoid_layer.hpp"
5
6    namespace caffe {
7
8    template <typename Dtype>
9    inline Dtype sigmoid(Dtype x) {
10     return 1. / (1. + exp(-x));
11   }
12
13   template <typename Dtype>
14   void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
15       const vector<Blob<Dtype>*>& top) {
16     const Dtype* bottom_data = bottom[0]->cpu_data();
17     Dtype* top_data = top[0]->mutable_cpu_data();
18     const int count = bottom[0]->count();
19     for (int i = 0; i < count; ++i) {
20       top_data[i] = sigmoid(bottom_data[i]);
21     }
22   }
23
24   template <typename Dtype>
25   void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
26       const vector<bool>& propagate_down,
27       const vector<Blob<Dtype>*>& bottom) {
28     if (propagate_down[0]) {
29       const Dtype* top_data = top[0]->cpu_data();
30       const Dtype* top_diff = top[0]->cpu_diff();
31       Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();
32       const int count = bottom[0]->count();
33       for (int i = 0; i < count; ++i) {
34         const Dtype sigmoid_x = top_data[i];
35         bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);
36       }
37     }
38   }
39
40   #ifdef CPU_ONLY
41   STUB_GPU(SigmoidLayer);
42   #endif
43
44   INSTANTIATE_CLASS(SigmoidLayer);
45
46
47   }  // namespace caffe
```

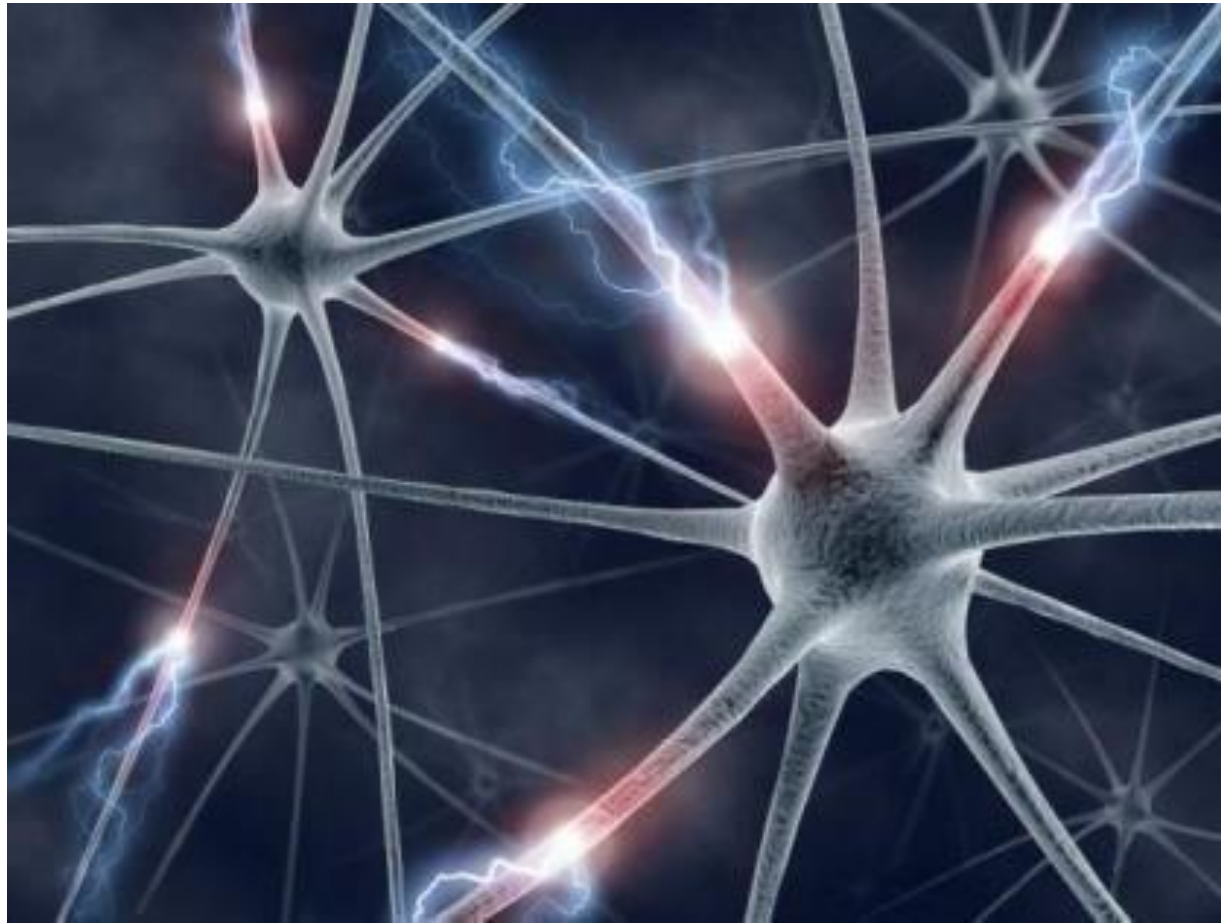forward()

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

backward()

$$(1 - \sigma(x))\,\sigma(x)$$ * top_diff  (chain rule)

# Backpropagation summary

- neural nets will be very large: impractical to write down gradient formula by hand for all parameters

- **backpropagation** = recursive application of the chain rule along a computational graph to compute the gradients of all inputs/parameters/intermediates

- implementations maintain a graph structure, where the nodes implement the **forward() / backward()** API

- **forward**: compute result of an operation and save any intermediates needed for gradient computation in memory

- **backward**: apply the chain rule to compute the gradient of the loss function with respect to the inputs

# Neural Networks

# Neural Networks

(**Before**) Linear score function:     $f = Wx$

# Neural Networks

(**Before**) Linear score function:    $f = Wx$

(**Now**) 2-layer Neural Network:    $f = W_2 \max(0, W_1 x)$

# Neural Networks

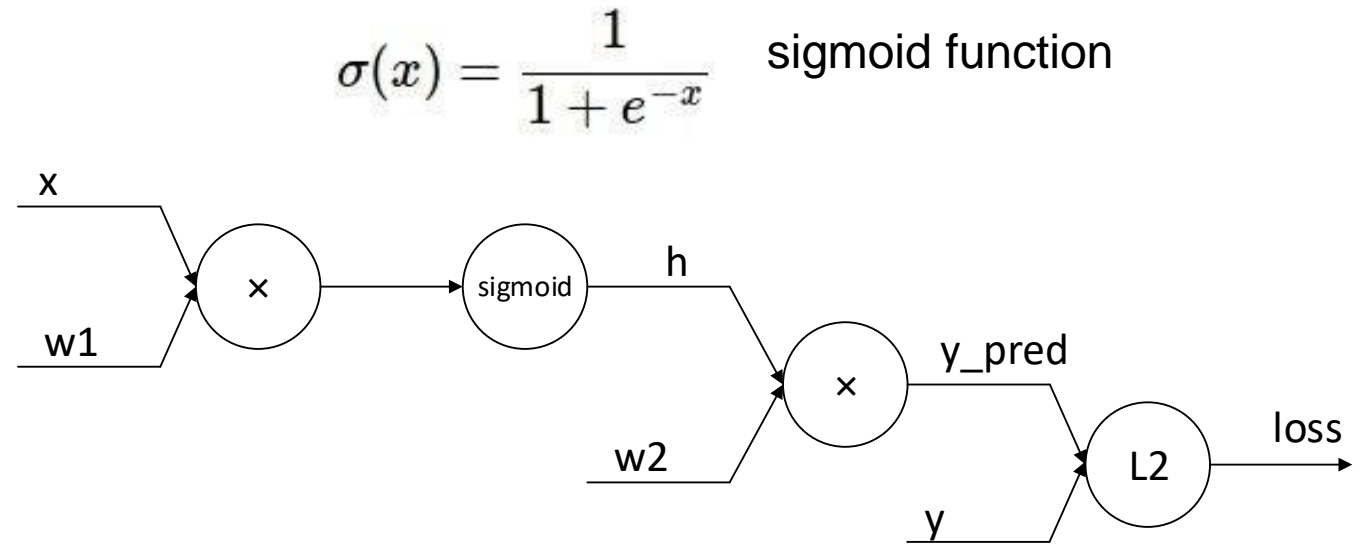(**Before**) Linear score function:    $f = Wx$

(**Now**) 2-layer Neural Network:    $f = W_2 \max(0, W_1 x)$



3072
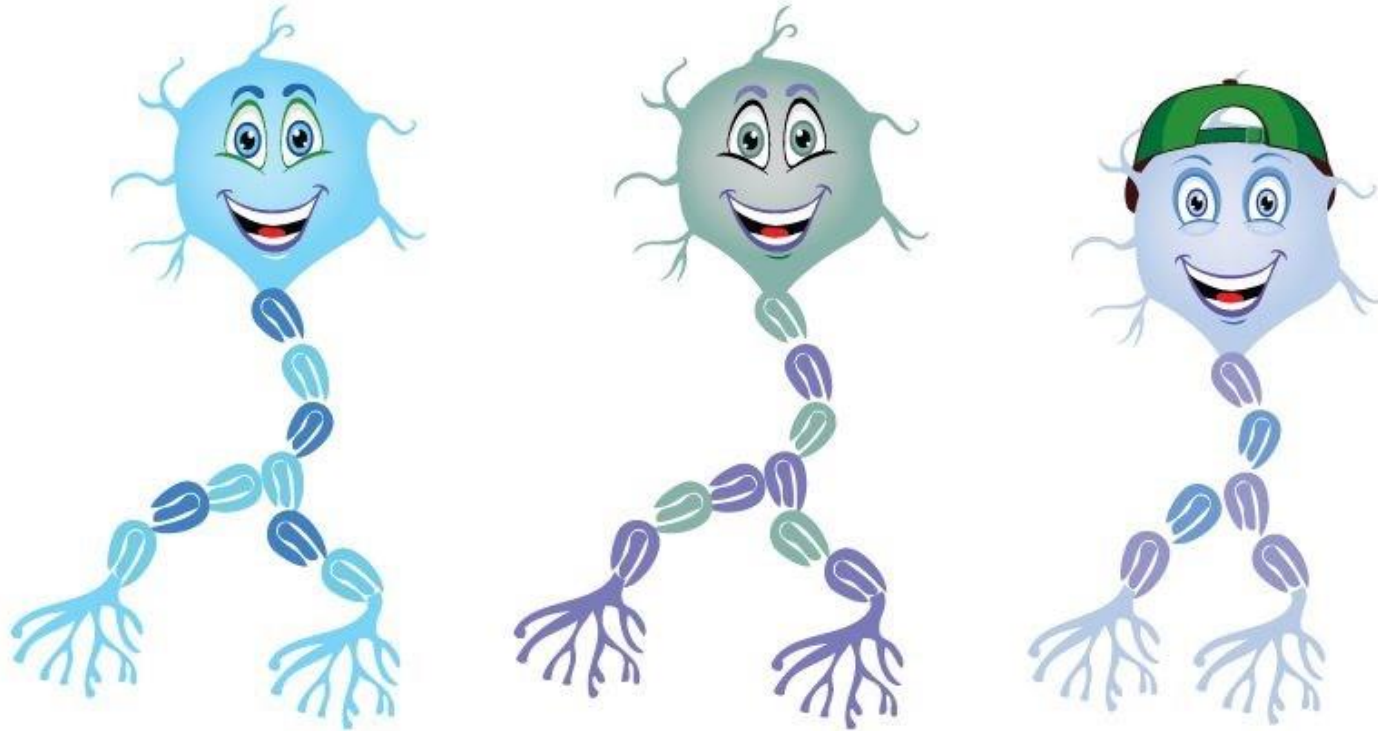
100

10

hidden layer

# Neural Networks

(**Before**) Linear score function:     $f = Wx$

(**Now**) 2-layer Neural Network:     $f = W_2 \max(0, W_1 x)$

# Neural Networks

(**Before**) Linear score function:     $f = Wx$

(**Now**) 2-layer Neural Network:     $f = W_2 \max(0, W_1 x)$

we can go deeper

3-layer Neural Network     $f = W_3 \max(0, W_2 \max(0, W_1 x))$

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```
1   import numpy as np
2   from numpy.random import randn
3
4   N, D_in, H, D_out = 64, 1000, 100, 10
5   x, y = randn(N, D_in), randn(N, D_out)
6   w1, w2 = randn(D_in, H), randn(H, D_out)
7
8   for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - y).sum()
12    print(t, loss)
13
14    grad_y_pred = 2.0 * (y_pred - y)
15    grad_w2 = h.T.dot(grad_y_pred)
16    grad_h = grad_y_pred.dot(w2.T)
17    grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19    w1 -= 1e-4 * grad_w1
20    w2 -= 1e-4 * grad_w2
```

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

# Artificial neuron

# Biological neuron

# Artificial neuron

# Artificial neuron
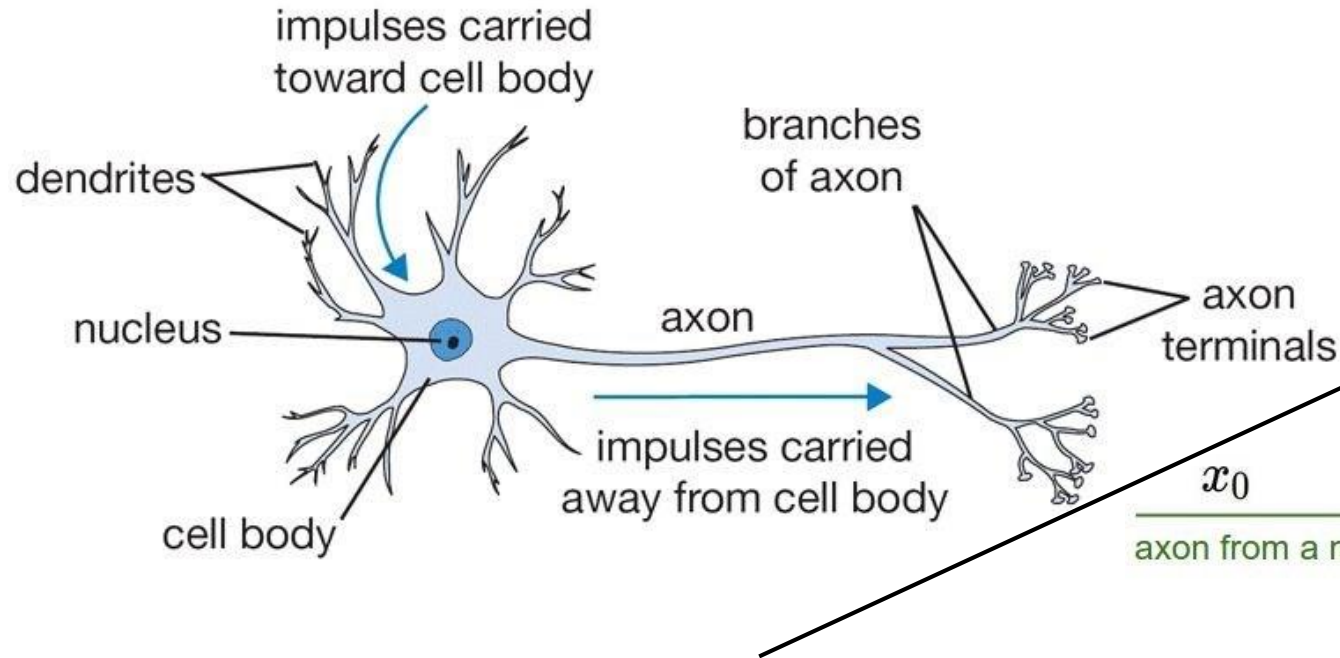


sigmoid activation function

$$\frac{1}{1+e^{-x}}$$

# Artificial neuron



impulses carried
toward cell body

dendrites

nucleus

cell body

branches
of axon

axon

impulses carried
away from cell body

axon
terminals

$x_0$  $w_0$

axon from a neuron    synapse

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$$\sum_i w_i x_i + b \quad f$$

$$f\left(\sum_i w_i x_i + b\right)$$

output axon

$w_2 x_2$

activation
function

```python
class Neuron:
    # ...
    def neuron_tick(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
        return firing_rate
```

# Активационные функции
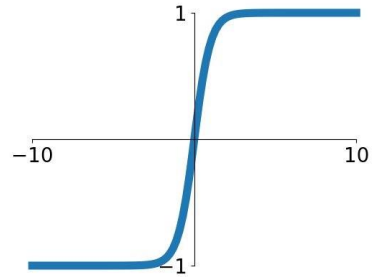
**Sigmoid**

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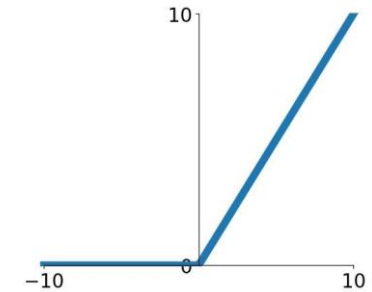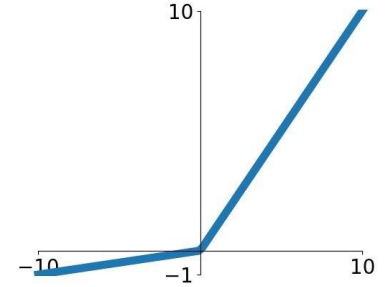
**tanh**

$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$

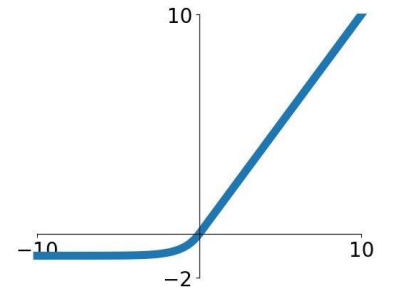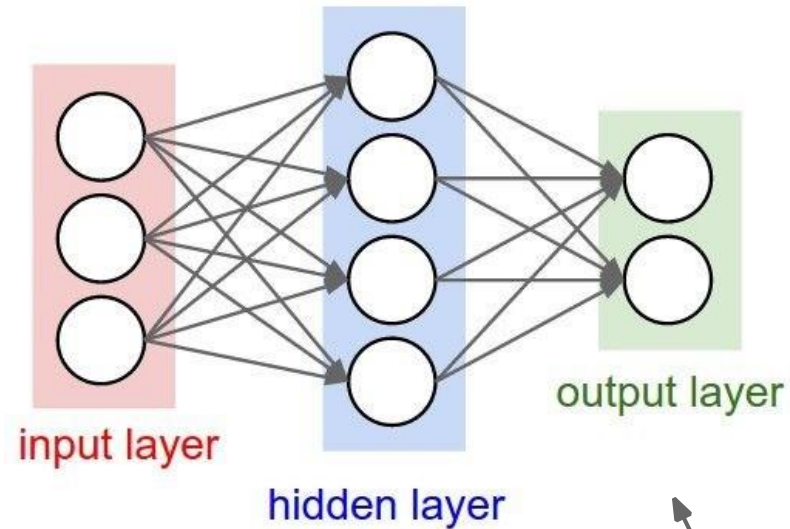**Leaky ReLU**

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

**Maxout neuron**

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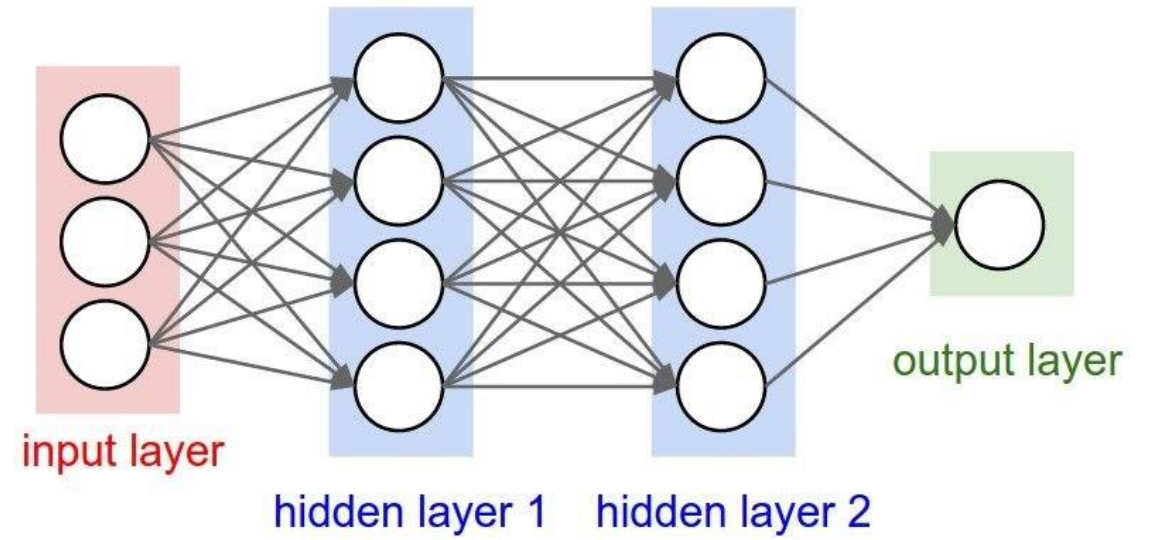
**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$
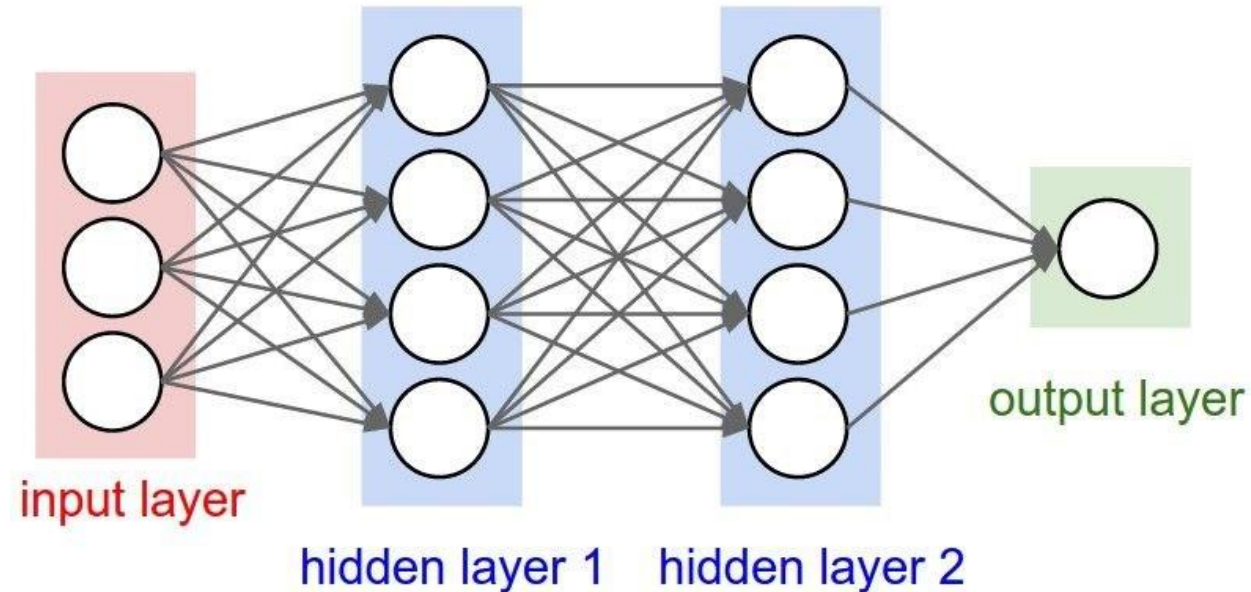
# Neural networks: fully-connected architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"
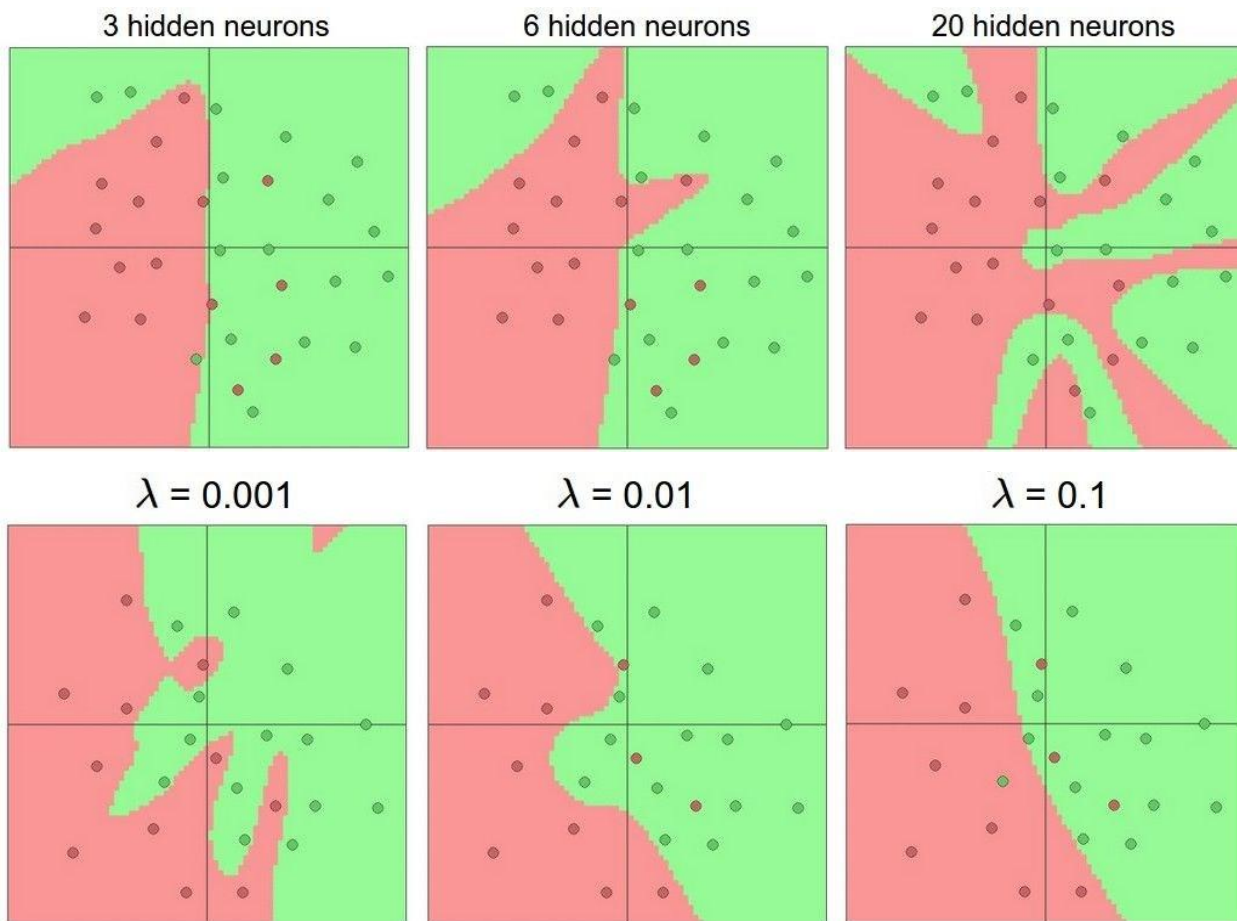
"3-layer Neural Net", or
"2-hidden-layer Neural Net"

**"Fully-connected" layers**

# Example feed-forward computation of a neural network



input layer

hidden layer 1   hidden layer 2

output layer

```
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

# Демо онлайн



Setting the number of layers and their sizes

Setting regularization

http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html

# В следующий раз

- Сверточные нейронные сети –
**Convolutional Neural Networks (CNN)**