# Lecture 4: Backpropagation and Neural Networks part 1

박사과정 김성빈 chengbinjin@inha.edu,
지도교수 김학일 교수 hikim@inha.ac.kr
인하대학교 컴퓨터비전 연구실

# Recall From Last Time…

Where we are…

$$s = f(x, W) = Wx$$

Score function

$$L_i = \sum_{j \neq y_i} \max\left(0, s_j - s_{y_i} + 1\right)$$

SVM loss

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Softmax loss

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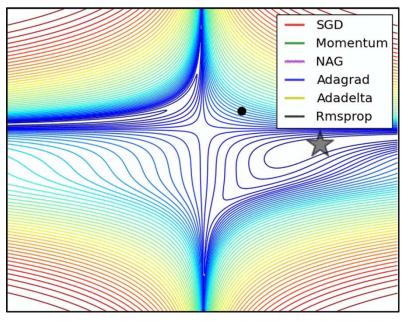
Data loss + Regularization loss

Want $\nabla_W L$ = …

gradient of Loss respect to W

# Recall From Last Time…

Optimization





```
# 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
```

(image credits to Alec Radford)

# Recall From Last Time…

Gradient Descent

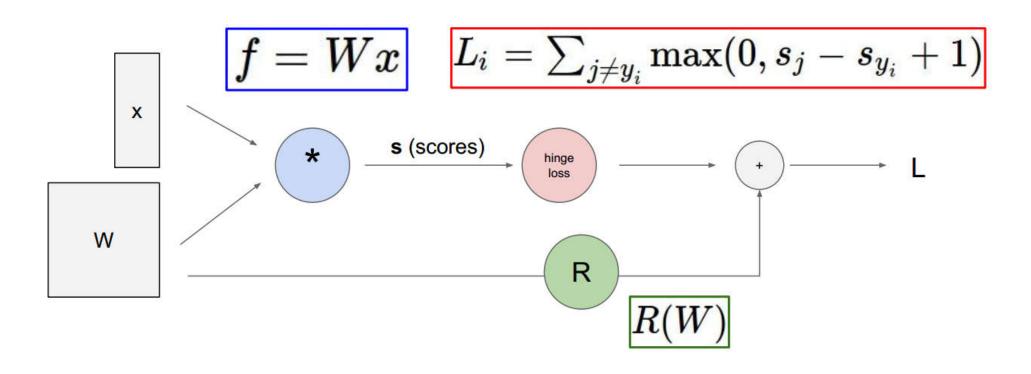$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

➢ **Numerical gradient**: approximate, slow, easy to write

➢ **Analytic gradient** (calculus): exact, fast, error-prone

➢ In practice: Derive analytic gradient, check your implementation with numerical gradient
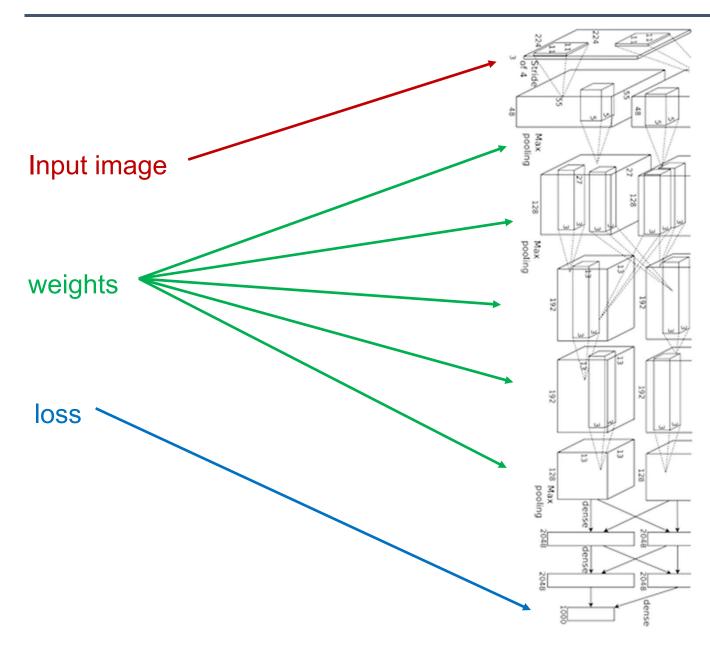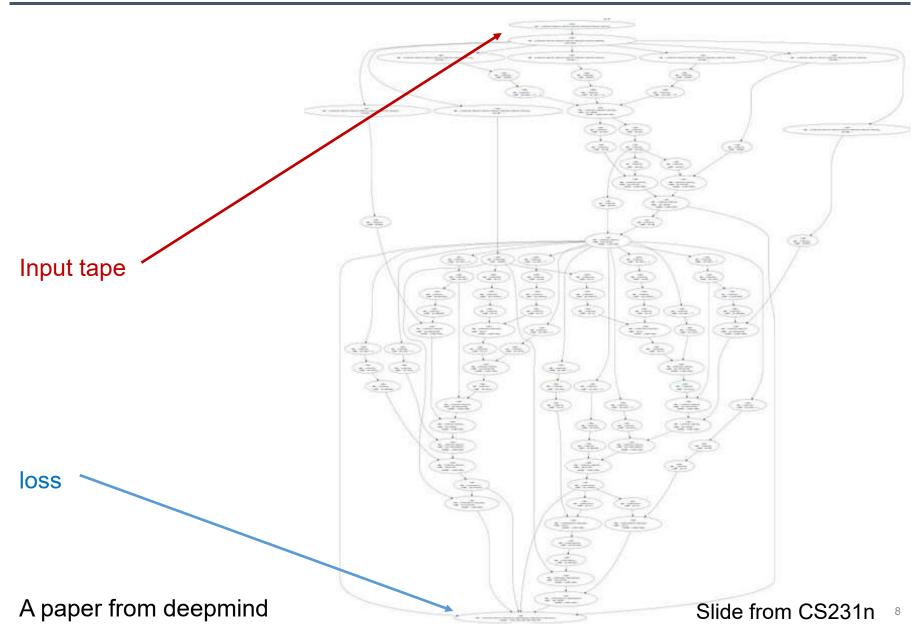
# Backpropagation

# Computational Graph

$$f = Wx \qquad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



Computation graph of the SVM loss

# Convolutional Network (AlexNet)

Input image

weights

loss

Slide from CS231n

# Neural Turing Machine



Input tape

loss

A paper from deepmind

# Neural Turing Machine

# Start from Simple

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

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

# Start from Simple

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

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

# Start from Simple

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

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

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

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



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

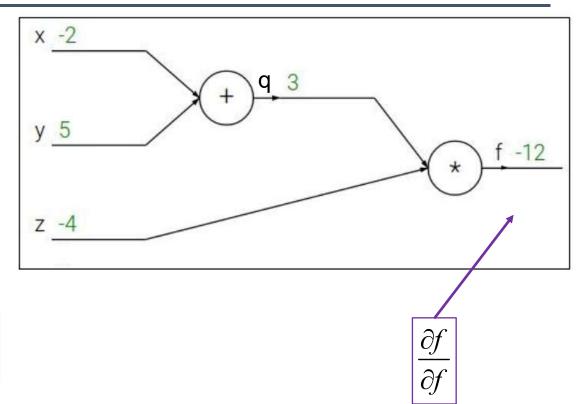Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

# Start from Simple

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

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



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

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

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

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

# Start from Simple

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

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

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

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



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

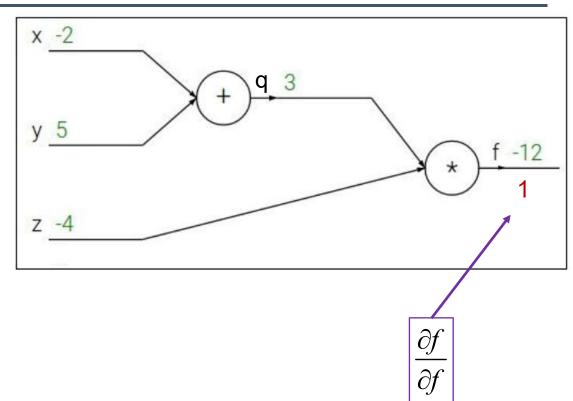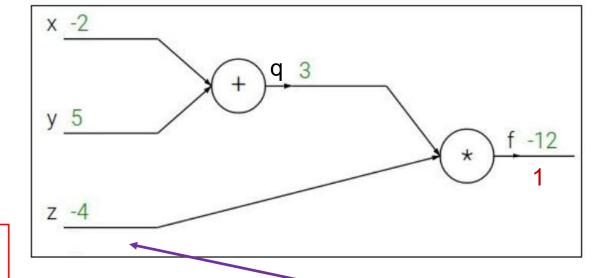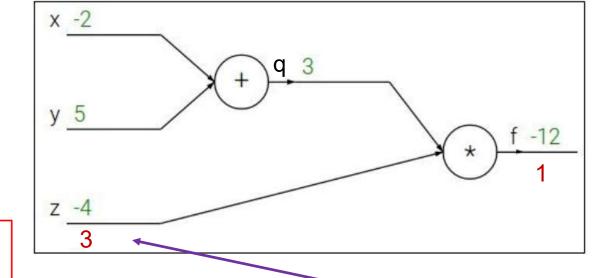Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Start from Simple

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

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

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

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



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

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

# Start from Simple

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

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

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

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



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

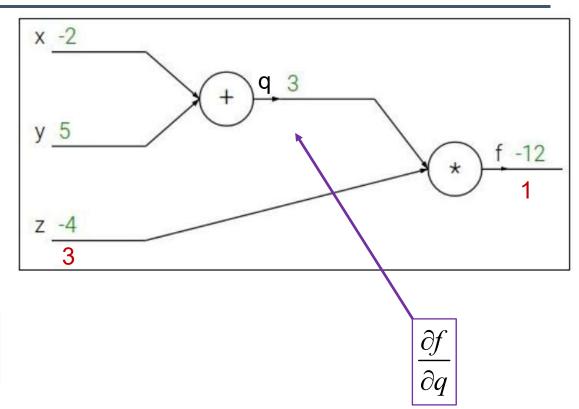Want: $\quad \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Start from Simple

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

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

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

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



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

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

# Start from Simple

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

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

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

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

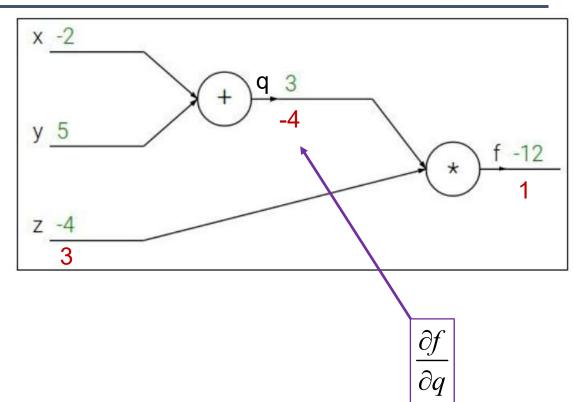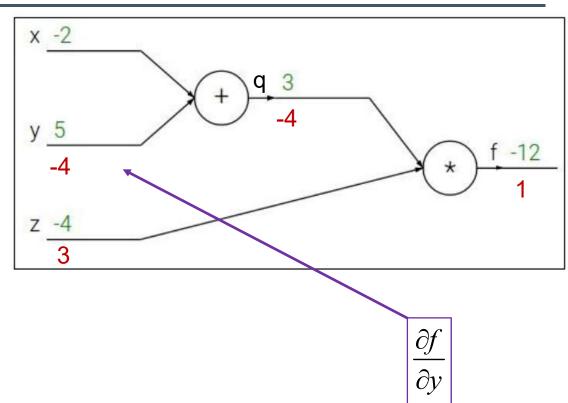Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

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

**Chain rule:**

$$\boxed{\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}}$$
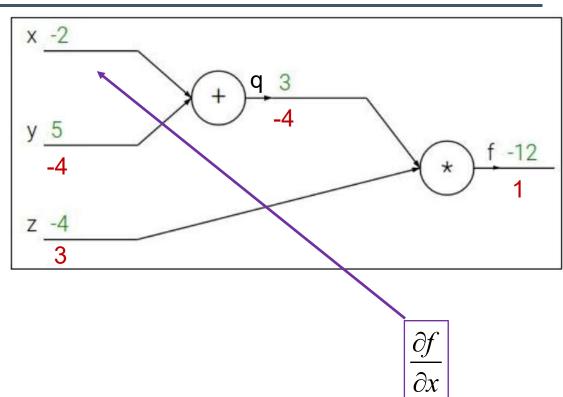
# Start from Simple

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

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

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

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

x  -2

q  3
+
-4

y  5
-4

f  -12
*
1

z  -4
3

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

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

# Start from Simple

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

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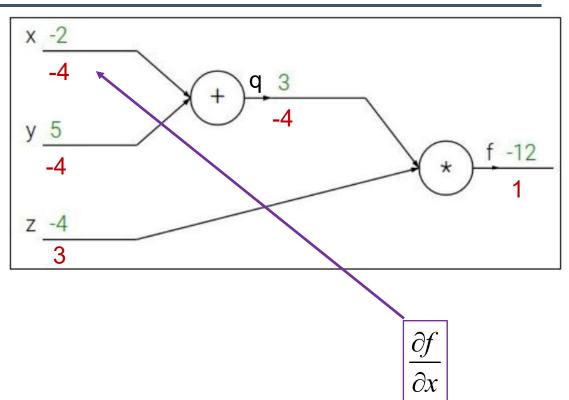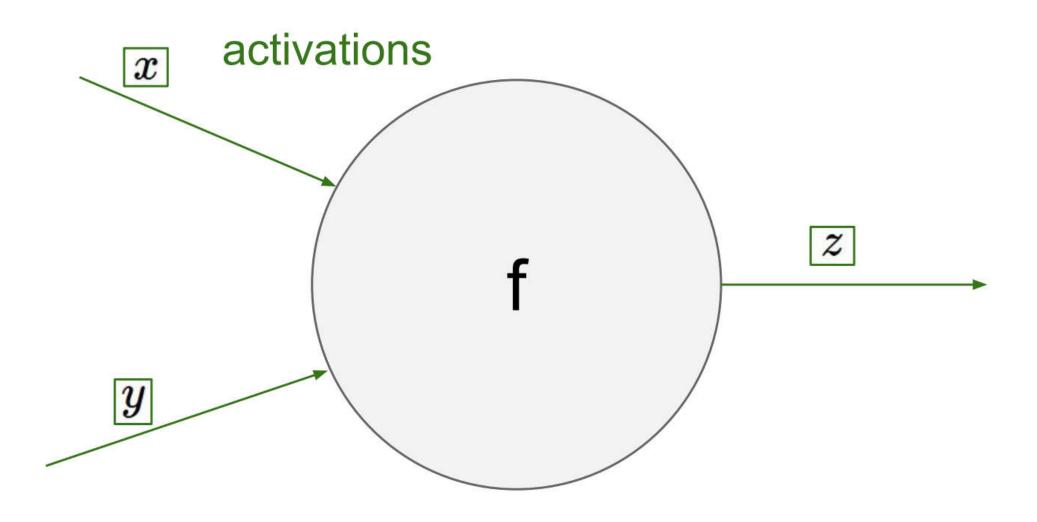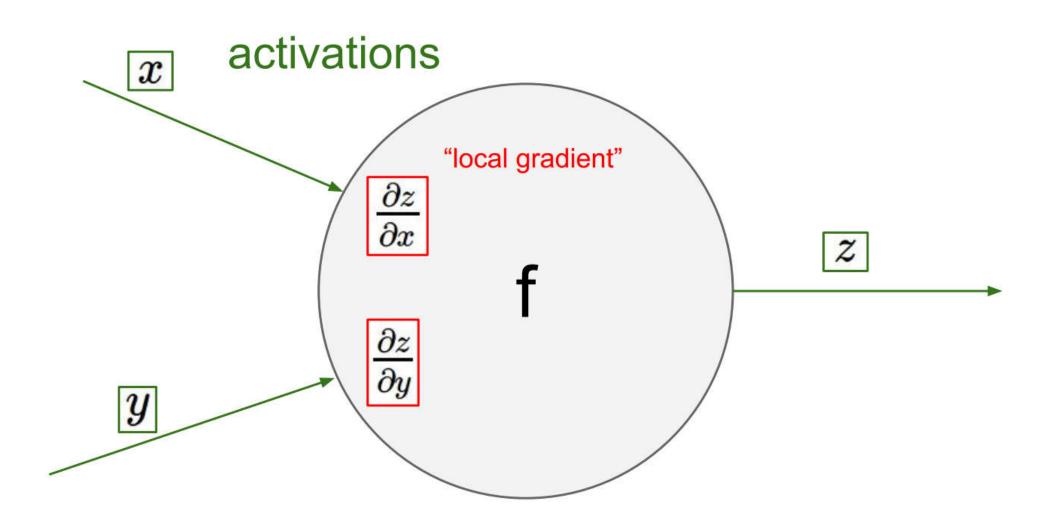


$\dfrac{\partial f}{\partial x}$

**Chain rule:**

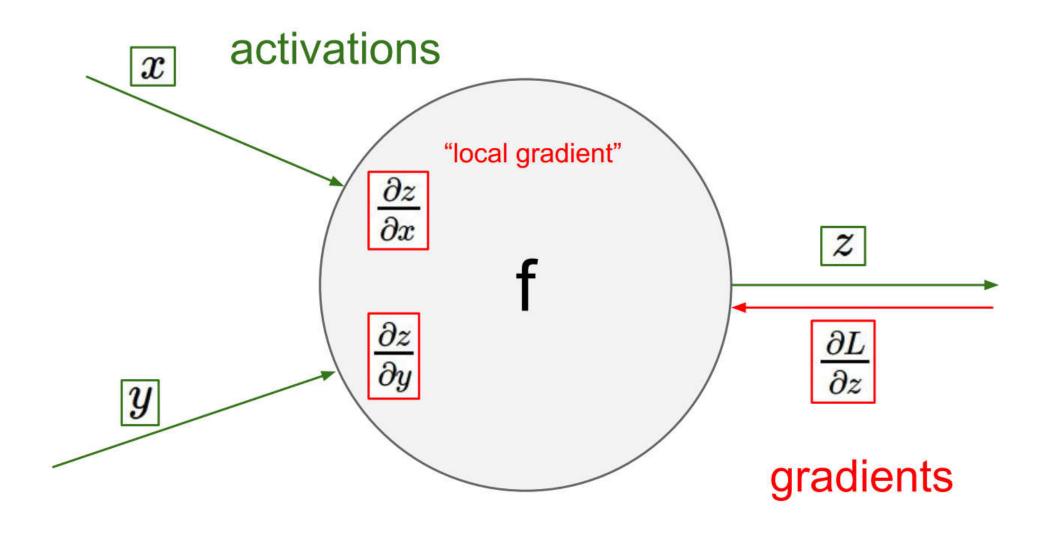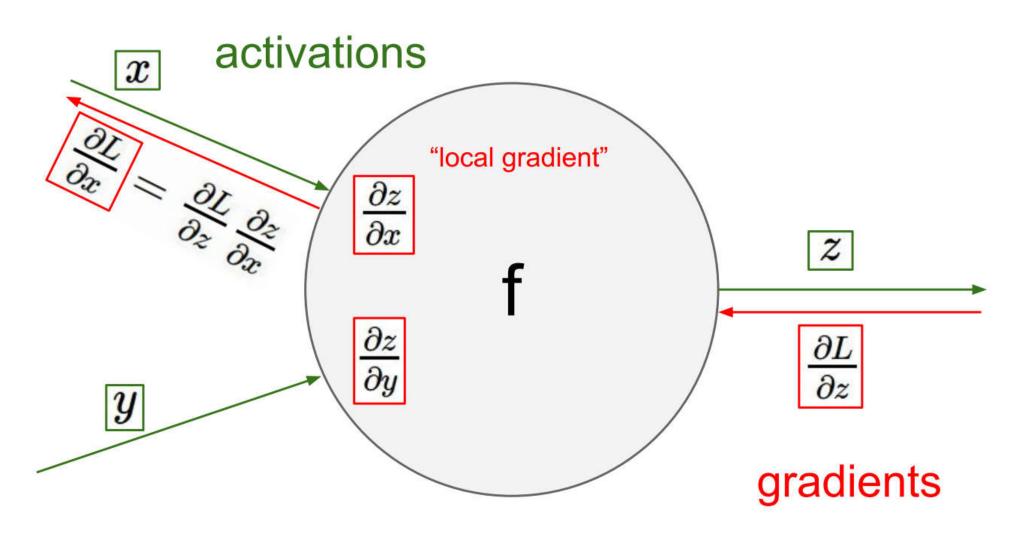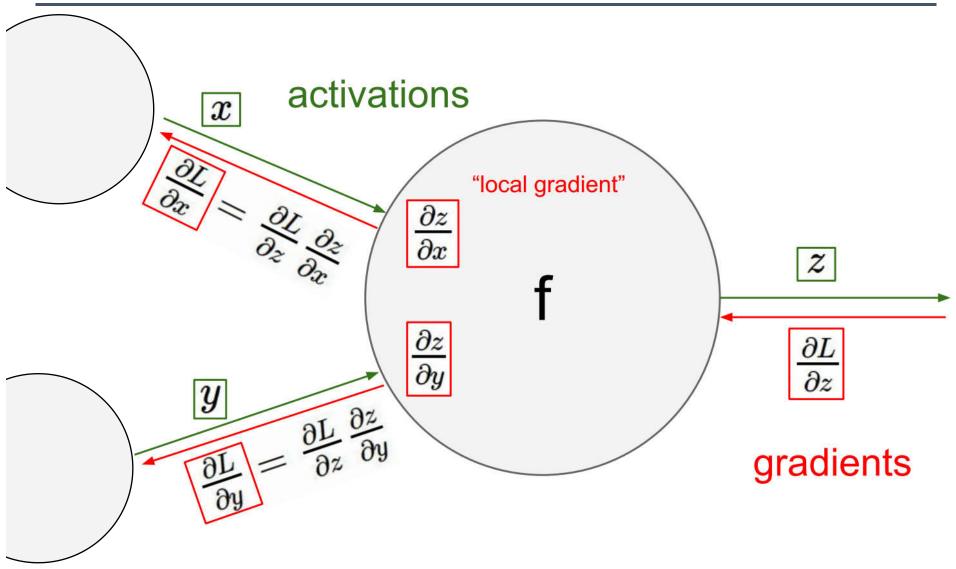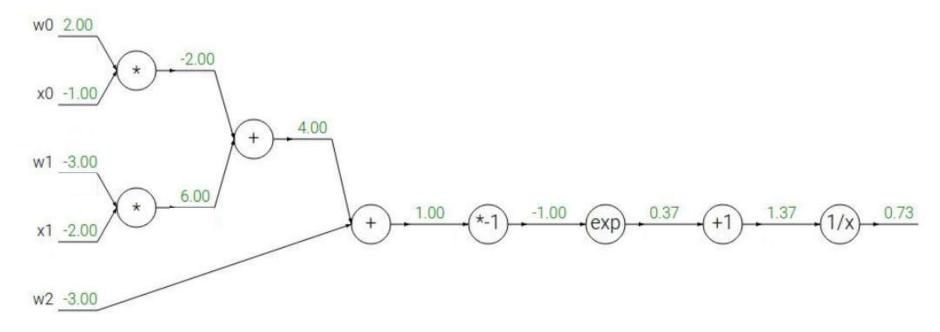$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q}\frac{\partial q}{\partial x}$$

activations

$x$

$y$

f

$z$

# Backpropagation

# Backpropagation

activations

$$x$$

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

"local gradient"

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

f

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

$$y$$

$$z$$

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

gradients

**Chain Rule**

Slide from CS231n 24

# Backpropagation



activations

"local gradient"

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

f

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

$x$

$y$

$z$

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

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

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

gradients

**Chain rule**

Slide from CS231n

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$
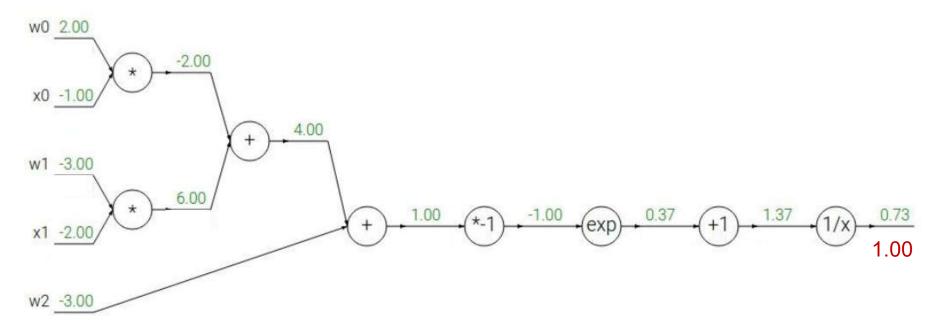


Computation graph of sigmoid neuron

# Another Example

**Sigmoid Neuron:** $\quad f(w,x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



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

Slide from CS231n

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00

-2.00

x0  -1.00

4.00

w1  -3.00

6.00

x1  -2.00          1.00        -1.00        exp        0.37        +1        1.37        1/x        0.73

1.00

w2  -3.00

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

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

Slide from CS231n  28

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00

x0  -1.00

w1  -3.00

x1  -2.00

w2  -3.00

-2.00

4.00

6.00

1.00   *-1   -1.00   exp   0.37   +1   1.37   1/x   0.73

1.00

> Gradient = Local gradient * upstream gradient

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

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

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

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

Slide from CS231n   29

# Another Example

**Sigmoid Neuron:** $f(w,x) = \dfrac{1}{1+e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00

x0  -1.00

-2.00

w1  -3.00

4.00

6.00

x1  -2.00

w2  -3.00

1.00    -1.00    0.37    1.37    0.73

*-1    exp    +1    1/x

-0.53    1.00

➢ Gradient = Local gradient * upstream gradient

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

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

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

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

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

Slide from CS231n    30

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

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

-0.53

1.00

w2  -3.00

➢ Gradient = Local gradient * upstream gradient

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

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

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

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

Slide from CS231n  31

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



w0  2.00

x0  -1.00

        -2.00

w1  -3.00

x1  -2.00

        6.00

        4.00

w2  -3.00

        1.00        -1.00        0.37        1.37        0.73

        -0.53        -0.53        1.00

➢ Gradient = Local gradient * upstream gradient

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

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

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

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

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

Slide from CS231n  32

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



➢ Gradient = Local gradient * upstream gradient

$$\left(e^{-1}\right) * (-0.53) = -0.20$$

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

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

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

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

Slide from CS231n  33

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



➤ Gradient = Local gradient * upstream gradient

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
-0.20
-0.53
-0.53
1.00

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

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

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

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

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



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

0.20

-0.20

-0.53

-0.53

1.00

➢ Gradient = Local gradient * upstream gradient

$(-1)*(-0.20) = 0.20$

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

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

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

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

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



> ➤ Gradient = Local gradient * upstream gradient

$(1)*(0.20) = 0.20$

$(1)*(0.20) = 0.20$

| | | |
|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\dfrac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\dfrac{df}{dx} = a$ |

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

Slide from CS231n   36

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00

x0  -1.00

-2.00

4.00

0.20

w1  -3.00

x1  -2.00

6.00

0.20

w2  -3.00

0.20

> Gradient = Local gradient * upstream gradient

+ → 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

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

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

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

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

Slide from CS231n

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00

-2.00
0.20

x0  -1.00

4.00
0.20

w1  -3.00

6.00
0.20

x1  -2.00

0.20

w2  -3.00
0.20

➢ Gradient = Local gradient * upstream gradient

$(1) * (0.20) = 0.20$

$(1) * (0.20) = 0.20$

1.00
0.20

-1.00
-0.20

0.37
-0.53

1.37
-0.53

0.73
1.00

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

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

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

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

# Another Example

**Sigmoid Neuron:** $f(w,x) = \dfrac{1}{1+e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0 2.00

x0 -1.00

-2.00

0.20

➢ Gradient = Local gradient * upstream gradient

w1 -3.00

x1 -2.00

6.00

0.20

4.00

0.20

1.00

0.20

-1.00

-0.20

0.37

-0.53

1.37

-0.53

0.73

1.00

w2 -3.00

0.20

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

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

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

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

Slide from CS231n

# Another Example

**Sigmoid Neuron:** $f(w,x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



➢ Gradient = Local gradient * upstream gradient

$x_0$: $(2)*(0.20) = 0.40$

$w_0$: $(-1)*(0.20) = -0.20$

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

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

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

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

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



➢ Gradient = Local gradient * upstream gradient

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

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

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

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

# Another Example

**Sigmoid Neuron:** $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00
-0.20

x0  -1.00
0.40

w1  -3.00
-0.40

x1  -2.00
-0.60

w2  -3.00
0.20

-2.00
0.20

4.00
0.20

6.00
0.20

1.00
0.20

-1.00
-0.20

0.37
-0.53

1.37
-0.53

0.73
1.00

(*) (+) (+) (*-1) (exp) (+1) (1/x)

➢ Gradient = Local gradient * upstream gradient

$x_1$: (-3)*(0.20) = -0.60

$w_1$: (-2)*(0.20) = -0.40

$f(x) = e^x$ → $\dfrac{df}{dx} = e^x$

$f_a(x) = ax$ → $\dfrac{df}{dx} = a$

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

$f_c(x) = c + x$ → $\dfrac{df}{dx} = 1$

# Sigmoid Function

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

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

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



**Sigmoid gate**

# Sigmoid Function

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

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

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



**Sigmoid gate**

[(1-0.73) *(0.73)] * 1.0=0.2

# Patterns in Backward Flow

add gate: gradient distributor

max gate: gradient router

mul gate: gradient "switcher"



max(-1, 2) = 2

max(-1+δ, 2) = 2

max(-1, 2+δ) = 2+δ

**Question #1: what will happen in max gate when two inputs are the same?**

# Gradients Add at Branches



- You will see this kind of operation in Batch Normalization

# Forward/Backward API

# Implementation: Forward/Backward API



**Rough pseudo 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
```

Slide from CS231n

# Implementation: Forward/Backward API



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

**Chain rule**

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

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

(x, y, z are scalars)

# Implementation: Forward/Backward API



```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]
```

cache

(x, y, z are scalars)

# Example: Torch Layers

# Example: Torch Layers



=

# Example: Torch MulConstant

```lua
local MulConstant, parent = torch.class('nn.MulConstant', 'nn.Module')

function MulConstant:__init(constant_scalar,ip)
  parent.__init(self)
  assert(type(constant_scalar) == 'number', 'input is not scalar!')
  self.constant_scalar = constant_scalar

  -- default for inplace is false
  self.inplace = ip or false
  if (ip and type(ip) ~= 'boolean') then
    error('in-place flag must be boolean')
  end
end

function MulConstant:updateOutput(input)
  if self.inplace then
    input:mul(self.constant_scalar)
    self.output = input
  else
    self.output:resizeAs(input)
    self.output:copy(input)
    self.output:mul(self.constant_scalar)
  end
  return self.output
end

function MulConstant:updateGradInput(input, gradOutput)
  if self.gradInput then
    if self.inplace then
      gradOutput:mul(self.constant_scalar)
      self.gradInput = gradOutput
      -- restore previous input value
      input:div(self.constant_scalar)
    else
      self.gradInput:resizeAs(gradOutput)
      self.gradInput:copy(gradOutput)
      self.gradInput:mul(self.constant_scalar)
    end
    return self.gradInput
  end
end
```

$$f(X) = aX$$

Initialization()

Forward()

Backward()

# Example: Caffe Layers
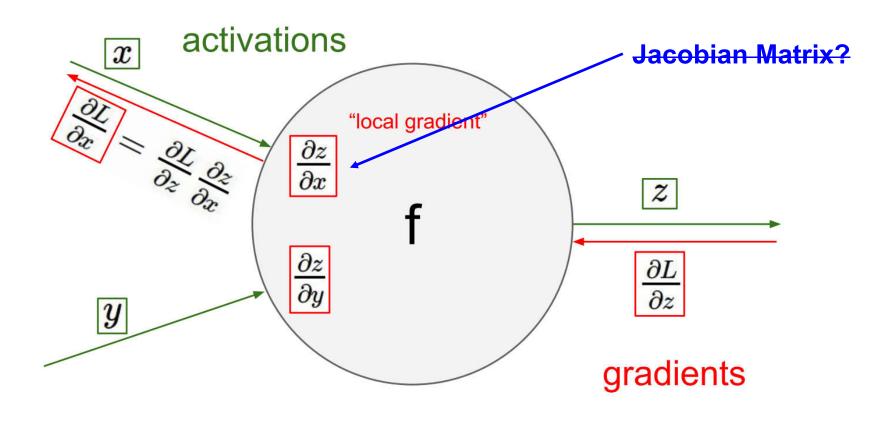
# Caffe Sigmoid Layer

```
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_cpu

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

Backward_cpu

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

Slide from CS231n  55

# Gradients for Vectorized Code

- Stage your forward/backward computation!

E.g. for the SVM:

```
# receive W (weights), X (data)
# forward pass (we have 8 lines)
scores = #...
margins = #...
data_loss = #...
reg_loss = #...
loss = data_loss + reg_loss
# backward pass (we have 5 lines)
dmargins = # ... (optionally, we go direct to dscores)
dscores = #...
dW = #...
```

margins

$$f = Wx \qquad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

x

W

\* s (scores)

hinge loss

R

$R(W)$

L

Computation graph of the SVM loss

# Summary so Far

- **Neural nets will be very large:** <u>no hope of writing down</u> gradient formula by hand for all parameters

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

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

- **Forward:** <u>compute result</u> of an operation and <u>save any intermediates</u> needed for gradient computation in memory

- **Backward:** apply the <u>chain rule</u> to compute the gradient of the loss function with respect to the inputs.

# Neural Network

# Neural Network: without the Brain Stuff
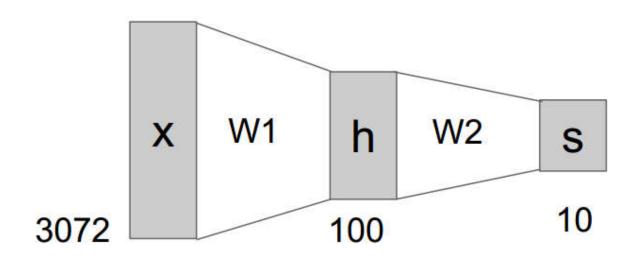
(Before) Linear score function: $f = Wx$

# Neural Network: without the Brain Stuff

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

(Now) 2-layer Neural Network $\quad f = W_2 \max(0, W_1 x)$ more complex mathematic expression of X

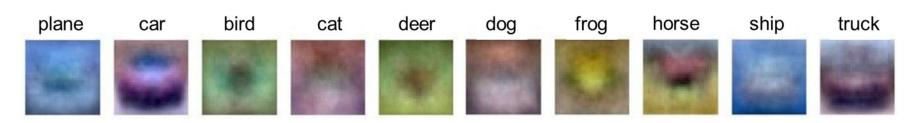# Neural Network: without the Brain Stuff

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

(Now) **2-layer Neural Network**  $\boxed{f = W_2\max(0, W_1 x)}$   more complex mathematic expression of X

➤ h is a hyperparameter, as big as possible fits in your computer



**mode**

| x | W1 | h | W2 | s |

3072        100        10

plane    car    bird    cat    deer    dog    frog    horse    ship    truck

# Neural Network: Without the Brain Stuff

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

(Now) 2-layer Neural Network $\quad f = W_2\max(0, W_1x)$ more complex mathematic expression of X

or **3-layer Neural Network** $\quad f = W_3\max(0, W_2\max(0, W_1x))$

# Python Code for 2-layer Neural Network

- Full implementation of training a 2-layer Neural Network needs ~11 lines:

```
01.  X = np.array([ [0,0,1],[0,1,1],[1,0,1],[1,1,1] ])
02.  y = np.array([[0,1,1,0]]).T
03.  syn0 = 2*np.random.random((3,4)) - 1
04.  syn1 = 2*np.random.random((4,1)) - 1
05.  for j in xrange(60000):
06.      l1 = 1/(1+np.exp(-(np.dot(X,syn0))))
07.      l2 = 1/(1+np.exp(-(np.dot(l1,syn1))))
08.      l2_delta = (y - l2)*(l2*(1-l2))
09.      l1_delta = l2_delta.dot(syn1.T) * (l1 * (1-l1))
10.      syn1 += l1.T.dot(l2_delta)
11.      syn0 += X.T.dot(l1_delta)
```

**Loss = (y – y_predict)**



➢ Sigmoid activation function

➢ Logistic regression loss

from @iamtrask, http://iamtrask.github.io/2015/07/12/basic-python-network/

# Assignment: Writing 2-Layer Net

- Stage your forward/backward computation!

```
# receive W1,W2,b1,b2 (weights/biases), X (data)
# forward pass:
h1 = #... function of X,W1,b1
scores = #... function of h1,W2,b2
loss = #... (several lines of code to evaluate Softmax loss)
# backward pass:
dscores = #...
dh1,dW2,db2 = #...
dW1,db1 = #...
```

# Neural Network Without Brain Stuff
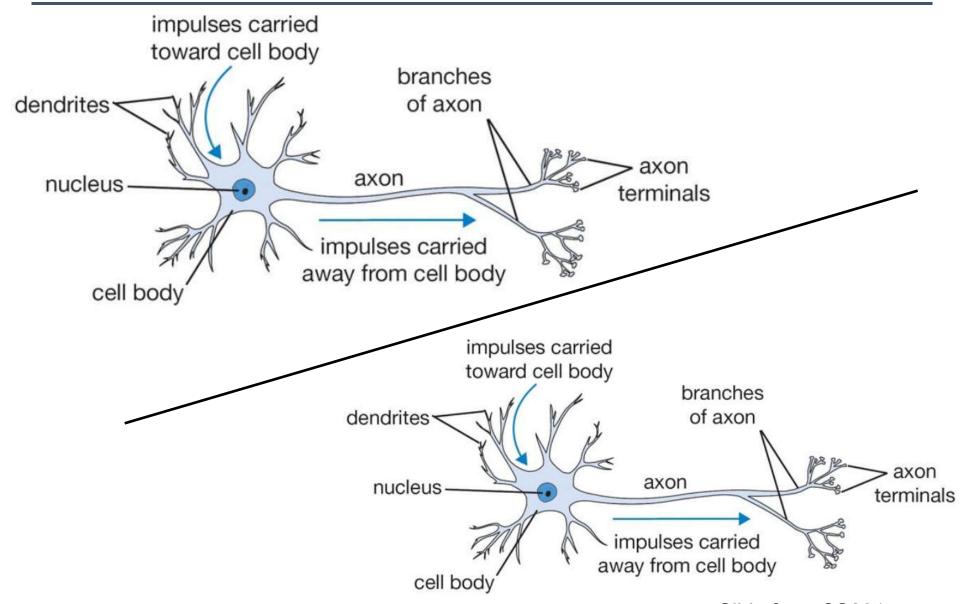
# Biological Neuron Without Brain Stuff

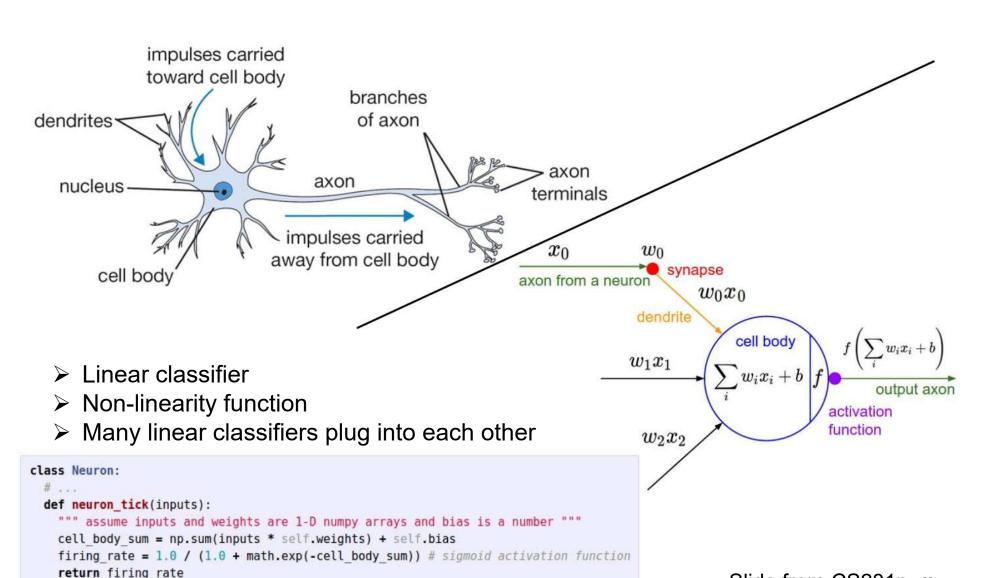# Biological Neuron Without Brain Stuff



**Sigmoid activation function**

# Biological Neuron Without Brain Stuff



> Linear classifier
> Non-linearity function
> Many linear classifiers plug into each other

```
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
```
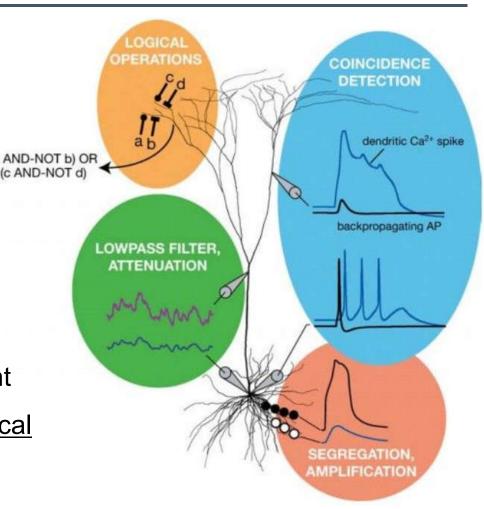
Slide from CS231n  69

# Biological Neuron

Be very careful with your Brain analogies:

**Biological Neurons:**

➤ Many different types of neurons

➤ **Dendrites** can perform complex non-linear computations

➤ **Synapses** are not a single weight but a complex non-linear dynamical system

➤ Rate code may not be adequate



LOGICAL OPERATIONS

(a AND-NOT b) OR (c AND-NOT d)

COINCIDENCE DETECTION

dendritic Ca²⁺ spike

backpropagating AP

LOWPASS FILTER, ATTENUATION

SEGREGATION, AMPLIFICATION

[Dendritic Computation. London and Hausser]
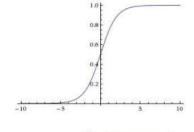
# Activation Functions

- Historically people use **Sigmoid and tanh**

- From 2012 **ReLu** became quite popular (default recommendation)

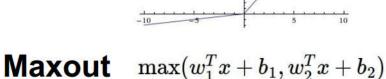- A few kind of hipster activation functions: **Leaky Relu**, **Maxout,** and **ELU**
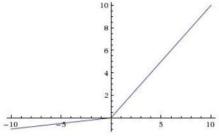
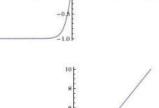**Question #2: why we need activation function?**

**Leaky ReLU**
$$\max(0.1x, x)$$

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

**tanh**  tanh(x)

**ReLU**  max(0,x)

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

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

Slide from CS231n  71

# Neural Networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"
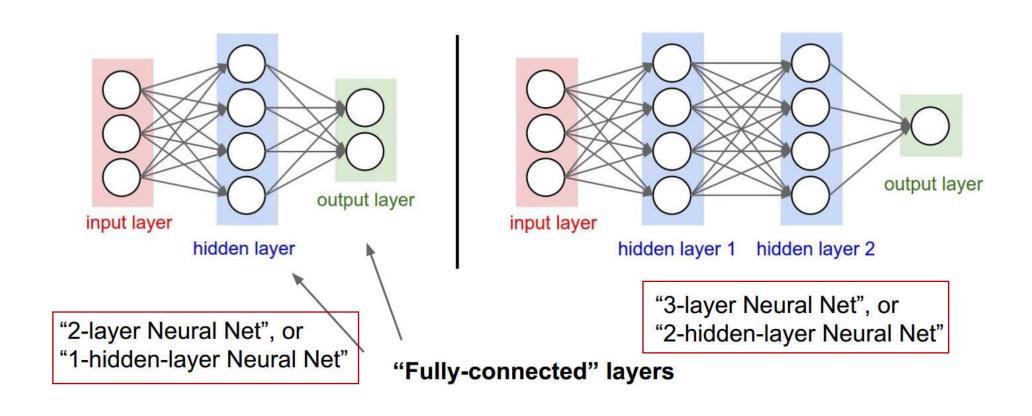
**"Fully-connected" layers**

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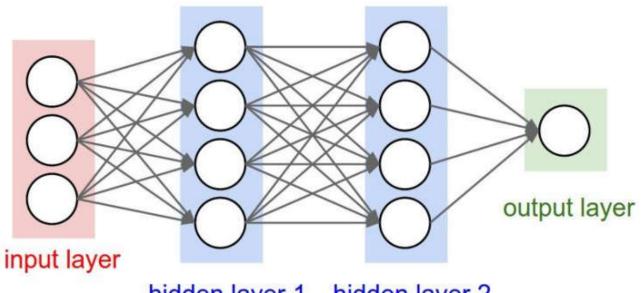"2-hidden-layer Neural Net"

# Example Feed-Forward Computation of a Neural Network

- Assume inputs and weights are 1-D numpy arrays and bias is a number

```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
```

# Example Feed-Forward Computation of a Neural Network

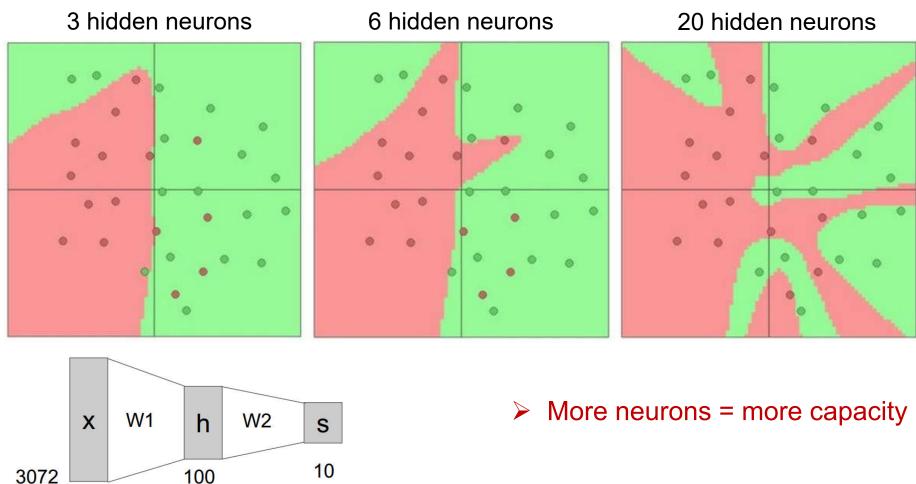- We can compute an entire set of neurons in a single hidden layer using **a single time matrix multiplication**



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)
```

Slide from CS231n

# Example of 2-Layer Neural Network
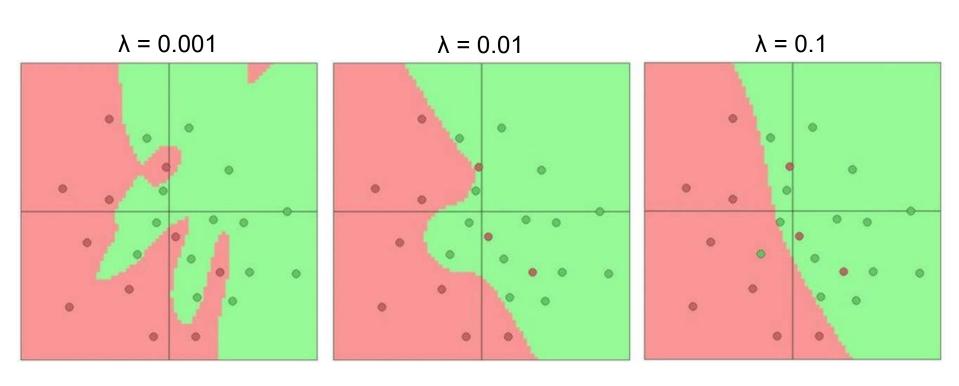
- Setting the **number of neurons** in the hidden layer

- h is **a hyperparameter**

3 hidden neurons        6 hidden neurons        20 hidden neurons



3072    x   W1   h   W2   s

100     10

➤ More neurons = more capacity

# Example of 2-Layer Neural Network

- Do not use size of neural network as a regularizer. Use **stronger regularization** instead

λ = 0.001              λ = 0.01              λ = 0.1



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

You can play with this demo over at ConvNetJs:
http://cs.Stanford.edu/people/karpathy/convnetjs/demo/classify2d.html

# Summary

- We arrange neurons into **fully-connected layers**

- The abstraction of a layer has the nice property that it allows us to use efficient **vectorized code** (e.g. **matrix multiplies**)

- Neural networks are **not really neural**

- Neural networks: **bigger = better** (<u>but might have to regularize more strongly</u>)

# Next Lecture

More than you ever wanted to know about

**Neural Networks** and **how to train them**.

# Thank you for your attention!