# Back-propagation in detail

Joseph Paul Cohen

# Goal for today

- Backpropagation (visualizing the chain rule)
- Intuition for applying gradient updates for arbitrary functions

# Back-Propagation

Described by Rumelhart, Hinton, and Williams in a 1986 paper

Method for efficiently calculating gradients in a multi-layer perceptron

Utilized by Yann LeCun in 90's to train first ConvNet

David
Rumelhart

Geoffrey
Hinton

## Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California, San Diego, La Jolla, California 92093, USA
† Department of Computer Science, Carnegie-Mellon University, Pittsburgh, Philadelphia 15213, USA
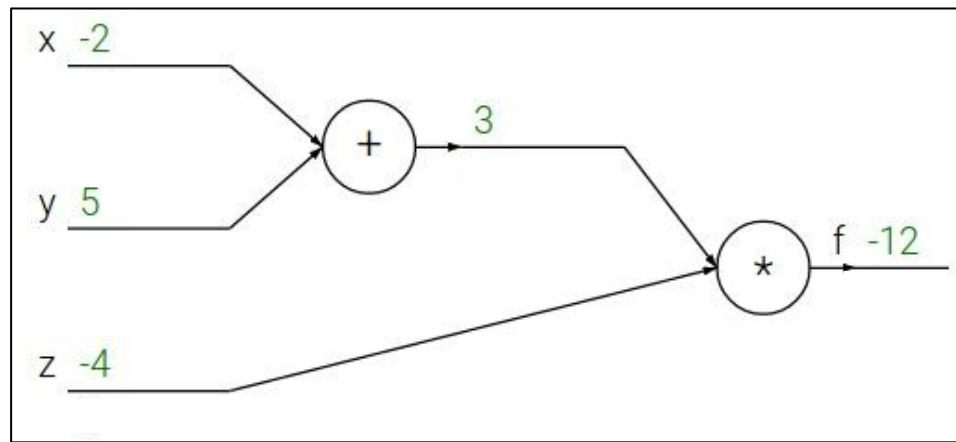
We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight

more difficult when we introduce hidden units whose actual or desired states are not specified by the task. (In perceptrons, there are 'feature analysers' between the input and output that are not true hidden units because their input connections are fixed by hand, so their states are completely determined by the input vector: they do not learn representations.) The learning procedure must decide under what circumstances the hidden units should be active in order to help achieve the desired input–output behaviour. This amounts to deciding what these units should represent. We demonstrate that a general purpose and relatively simple procedure is powerful enough to construct appropriate internal representations.

The simplest form of the learning procedure is for layered networks which have a layer of input units at the bottom; any number of intermediate layers; and a layer of output units at the top. Connections within a layer or from higher to lower layers are forbidden, but connections can skip intermediate layers. An input vector is presented to the network by setting

3

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

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

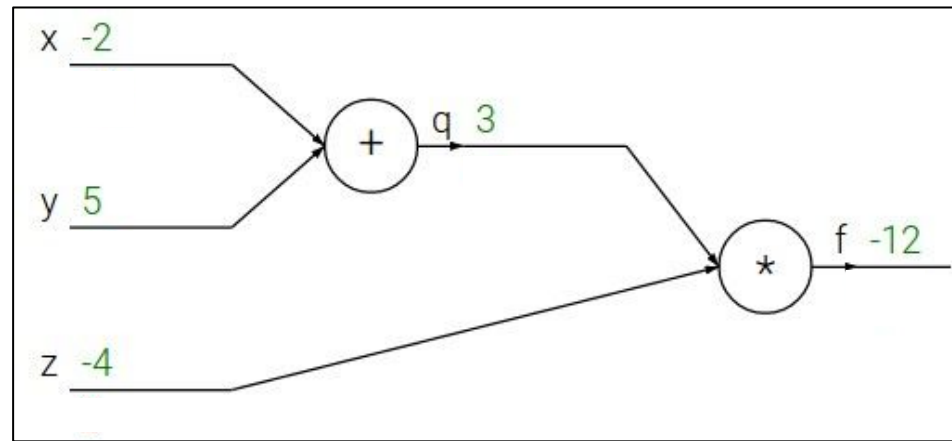$$f(x, y, z) = (x + y)z$$

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: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
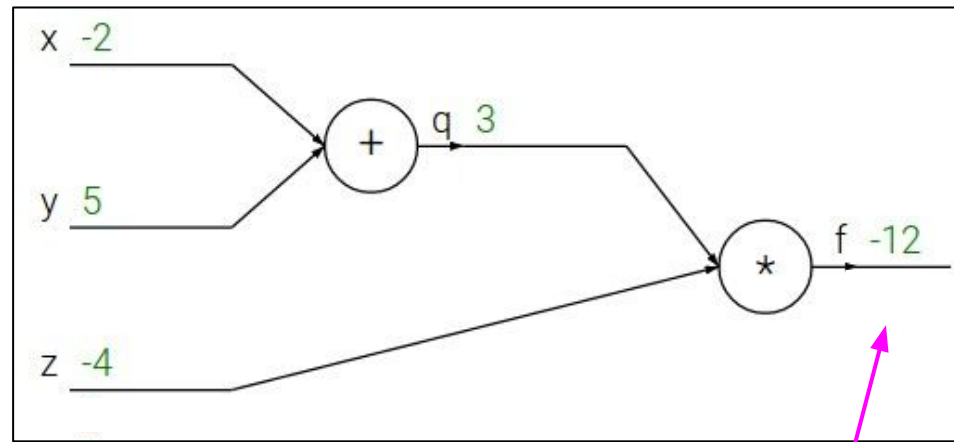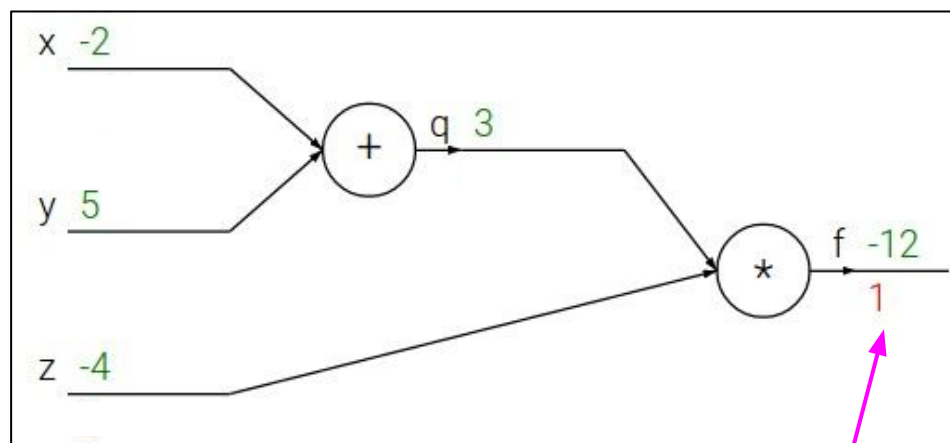
$$f(x, y, z) = (x + y)z$$

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

6

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

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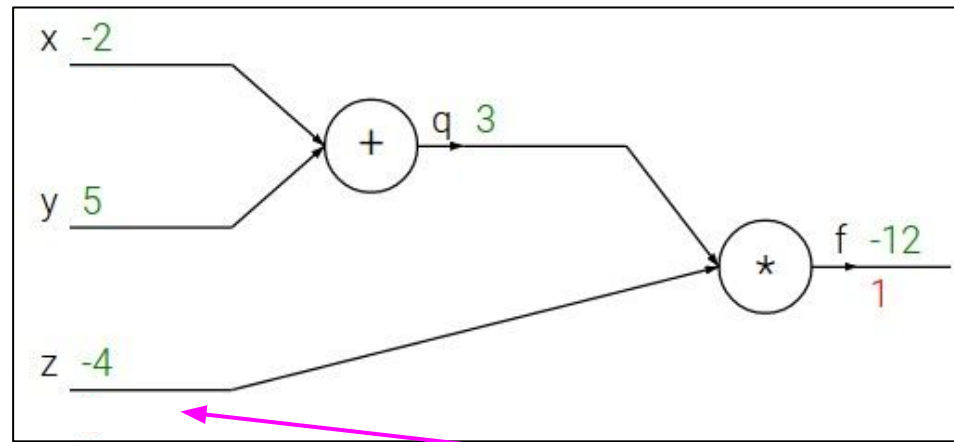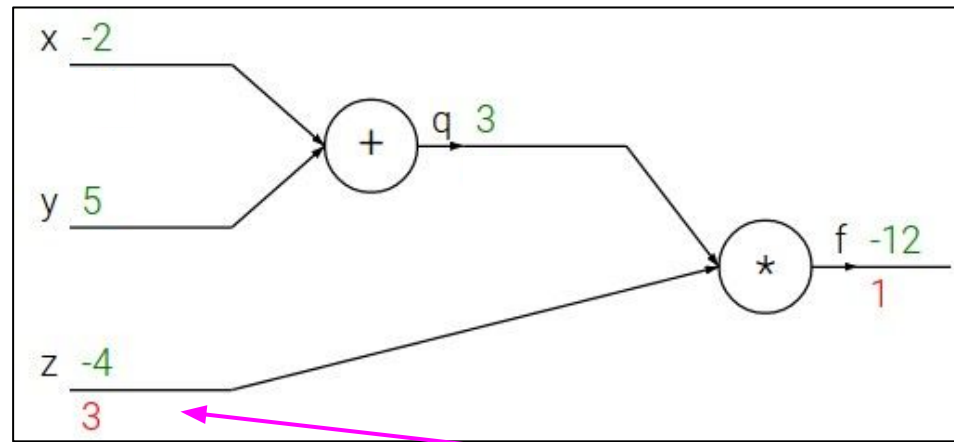


$$\frac{\partial f}{\partial f}$$

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

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

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

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

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

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

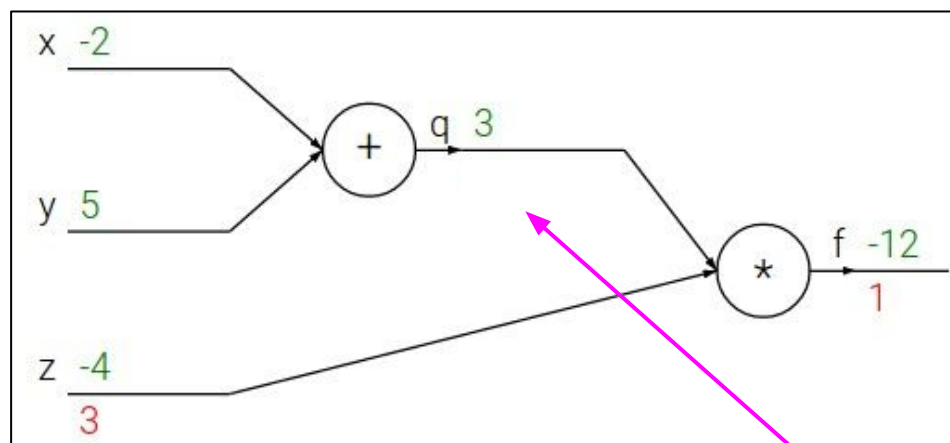Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
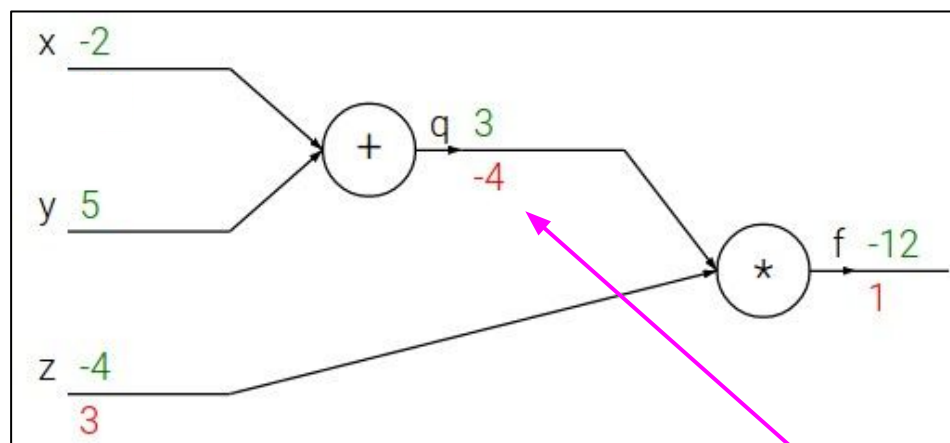


$$\frac{\partial f}{\partial z}$$

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

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

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

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



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

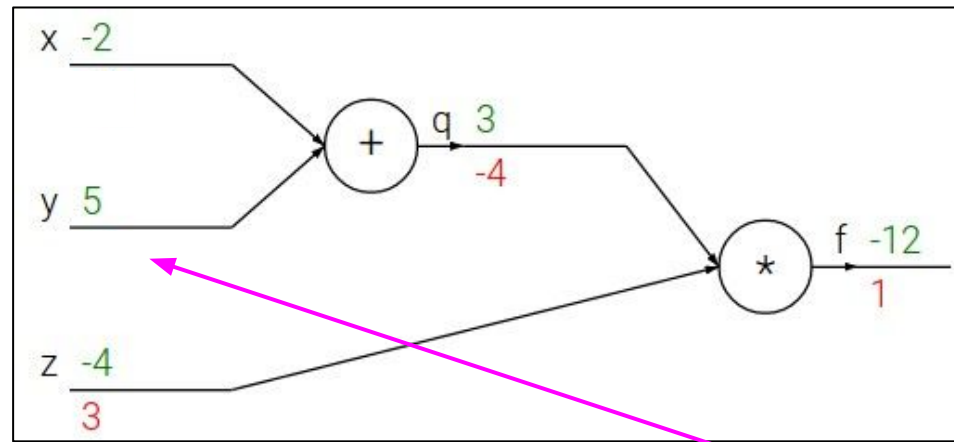Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

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

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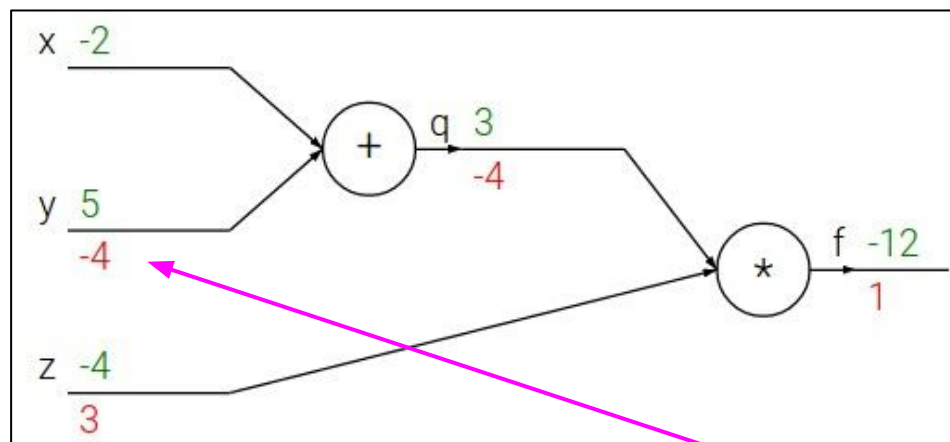


$$\frac{\partial f}{\partial q}$$

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

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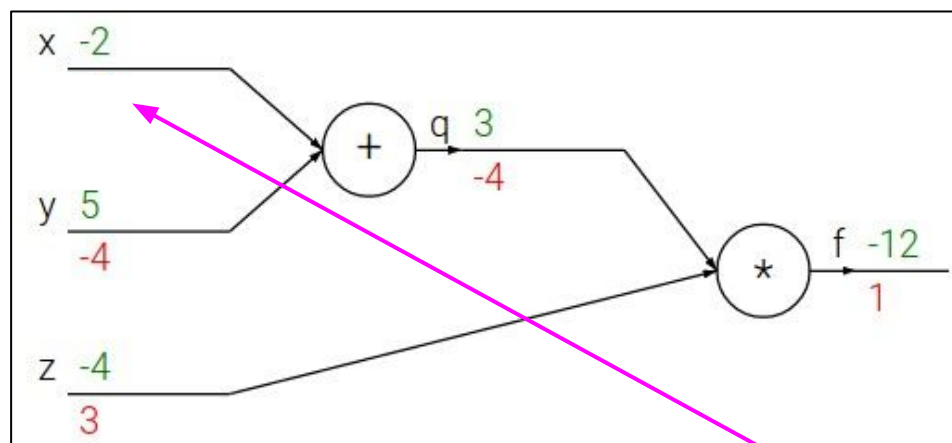


$$\frac{\partial f}{\partial y}$$

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

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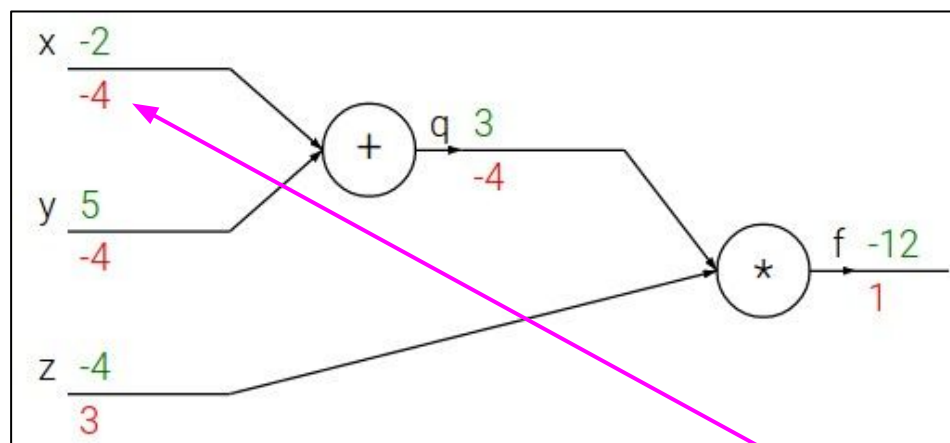


$$\frac{\partial f}{\partial y}$$

Chain rule:

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

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

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

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

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



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

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

14

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

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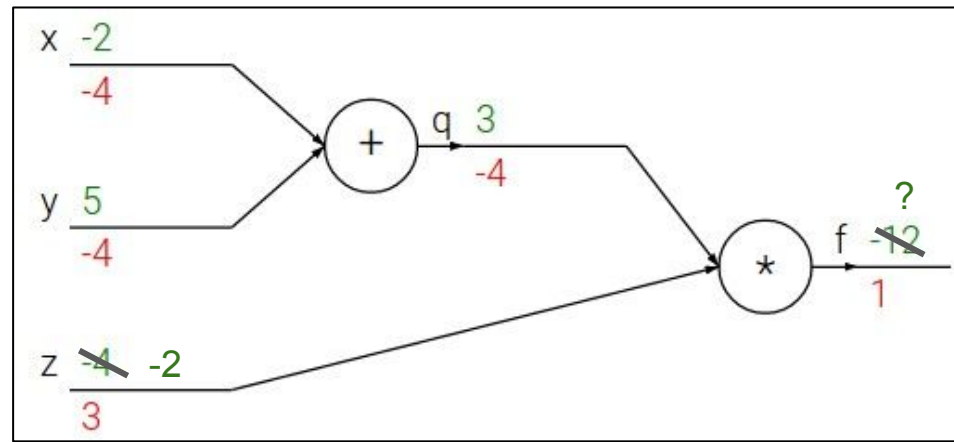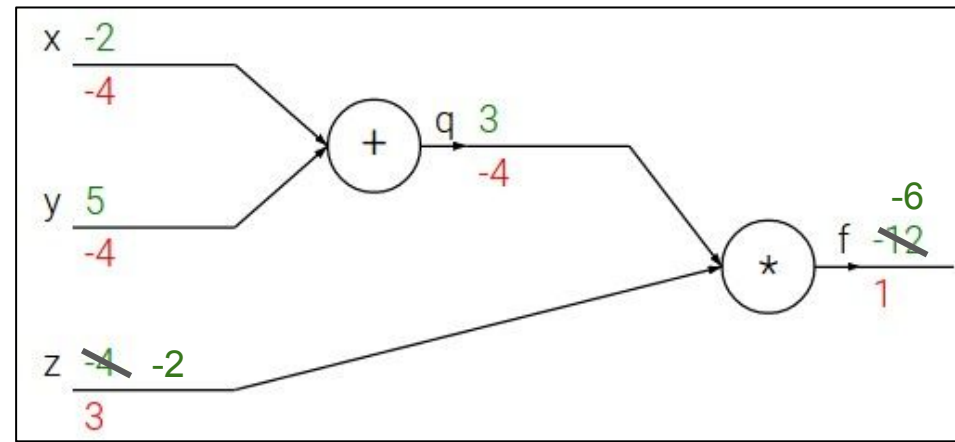


$$\frac{\partial f}{\partial x}$$

Chain rule:

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

15

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

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



x -2
-4

y 5
-4

+ q 3
-4

z -4 -2
3

* f -12
1

?

What happens to f when we add 2 to z?

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

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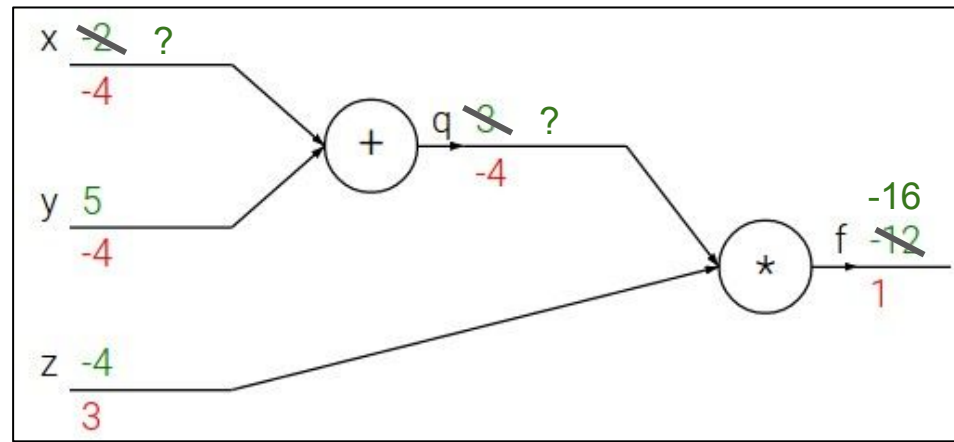


What happens to f when we add 2 to z?

z=z+2 causes f=f+(2*3)=f+6

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

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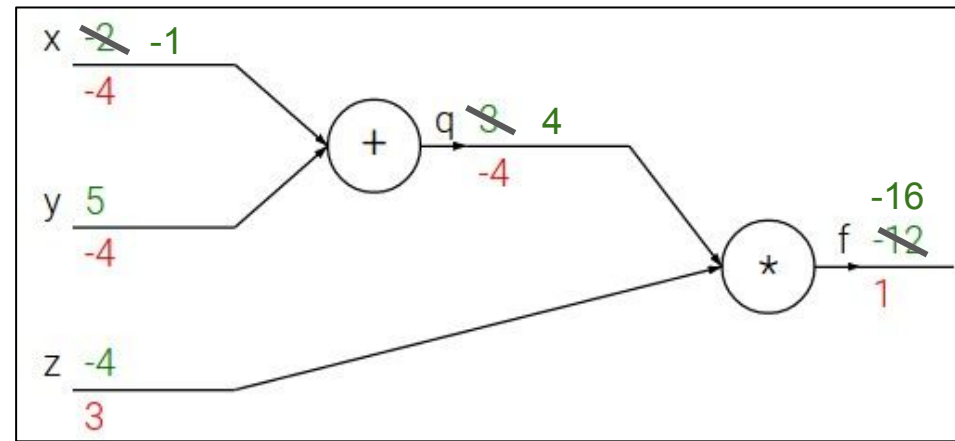


What happens to f when we add 2 to z?

z=z+2 causes f=f+(2*3)=f+6

How can we adjust x so that f is -16?

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

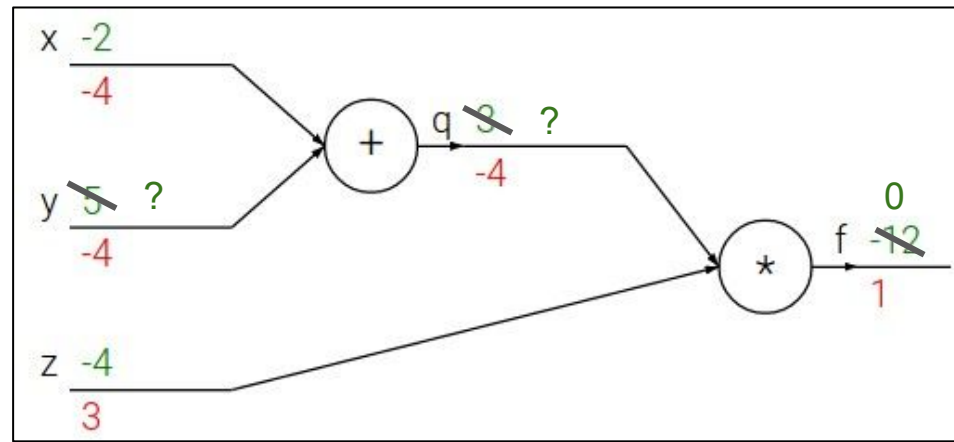e.g. x = -2, y = 5, z = -4



What happens to f when we add 2 to z?

z=z+2 causes f=f+(2*3)=f+6

How can we adjust x so that f is -16?

x=x+1 causes f=f+(1*-4)=f-4

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

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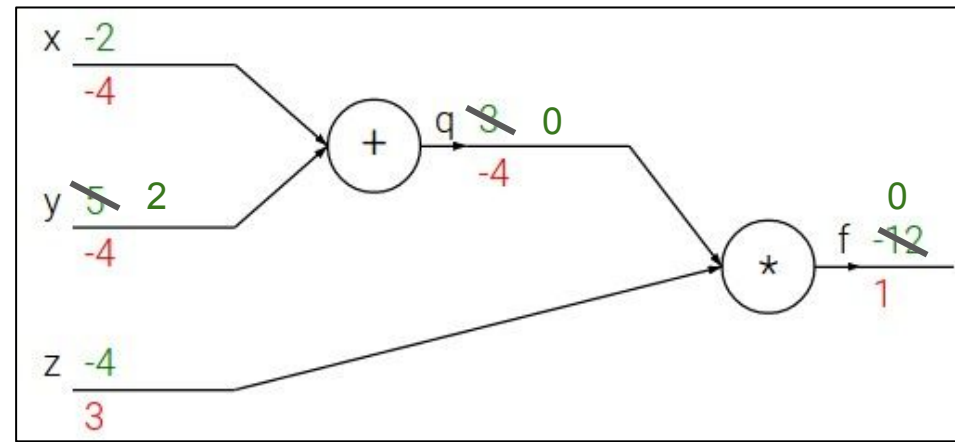


What happens to f when we add 2 to z?

z=z+2 causes f=f+(2*3)=f+6

How can we adjust x so that f is -16?

x=x+1 causes f=f+(1*-4)=f-4

How can we adjust y so that f is 0?

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

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



What happens to f when we add 2 to z?

z=z+2 causes f=f+(2*3)=f+6

How can we adjust x so that f is -16?
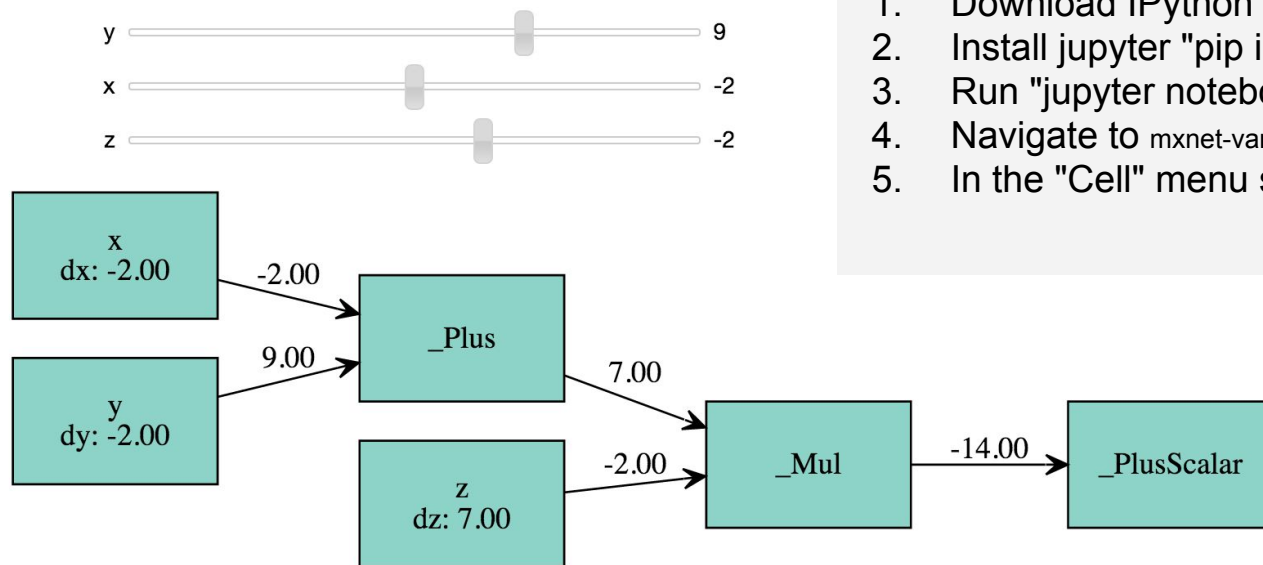
x=x+1 causes f=f+(1*-4)=f-4

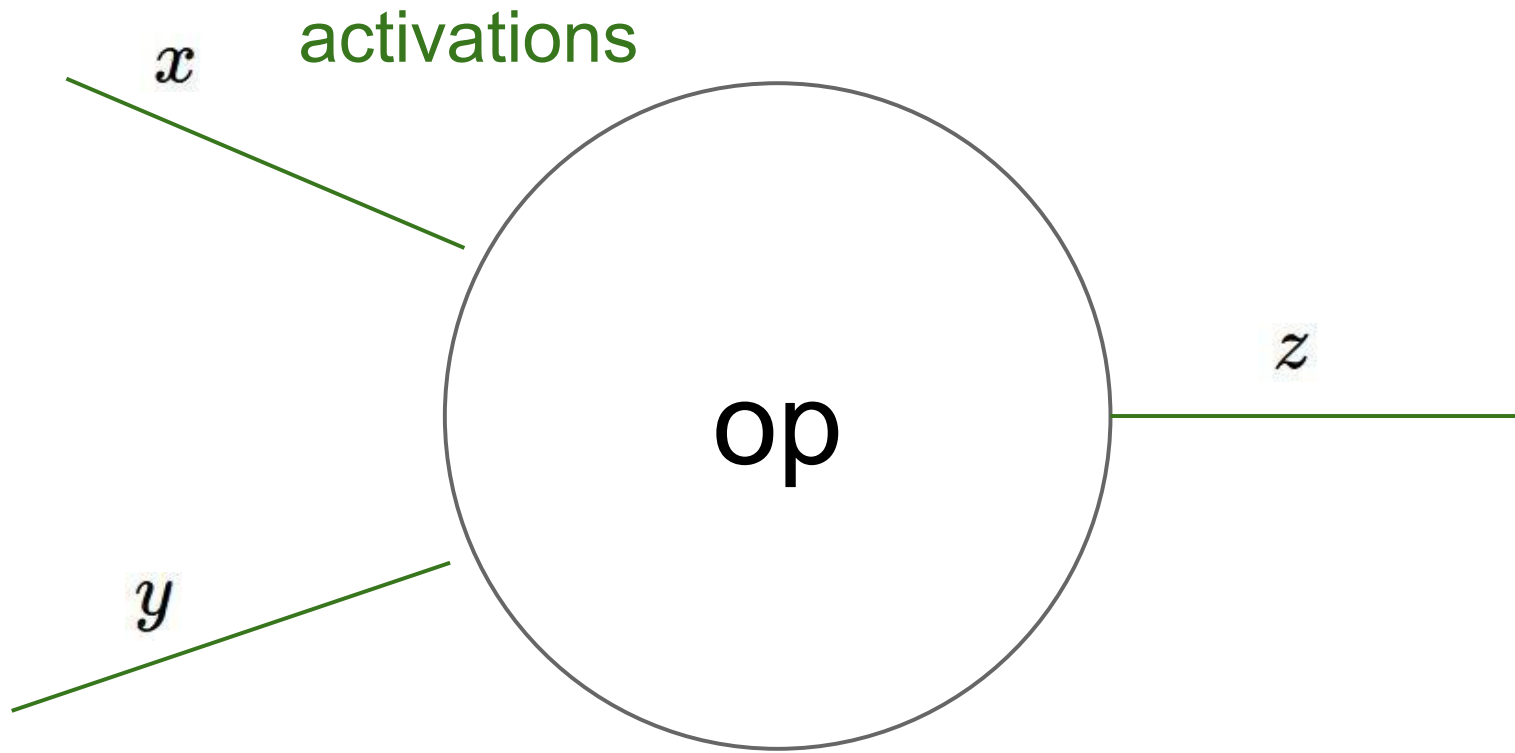How can we adjust y so that f is 0?

y=y-3 causes f=f+(-3*-4)=f+12

# Interactive demo


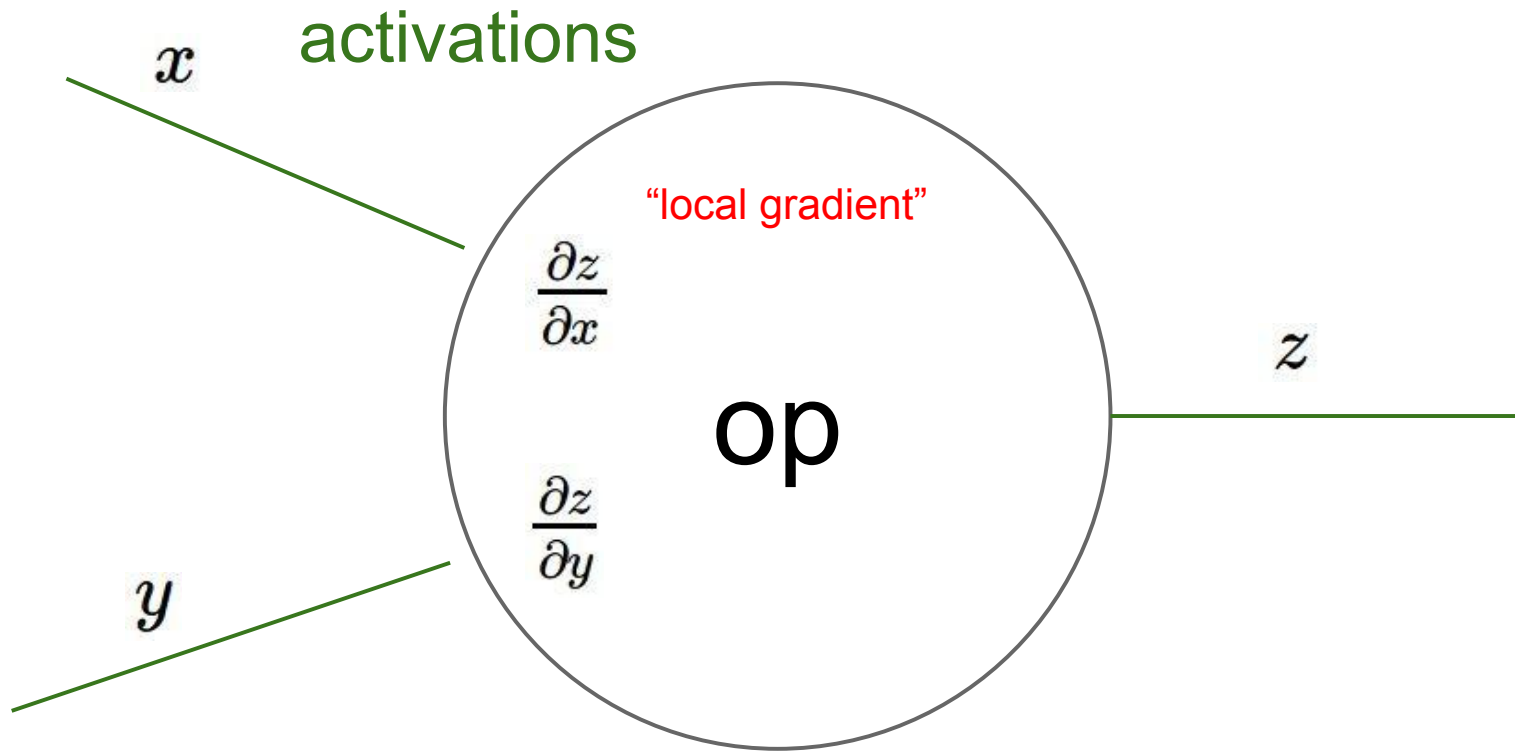
1. Download IPython notebook from github
2. Install jupyter "pip install jupyter"
3. Run "jupyter notebook"
4. Navigate to mxnet-vary-inputs-slideexamples.ipynb
5. In the "Cell" menu select "Run All"

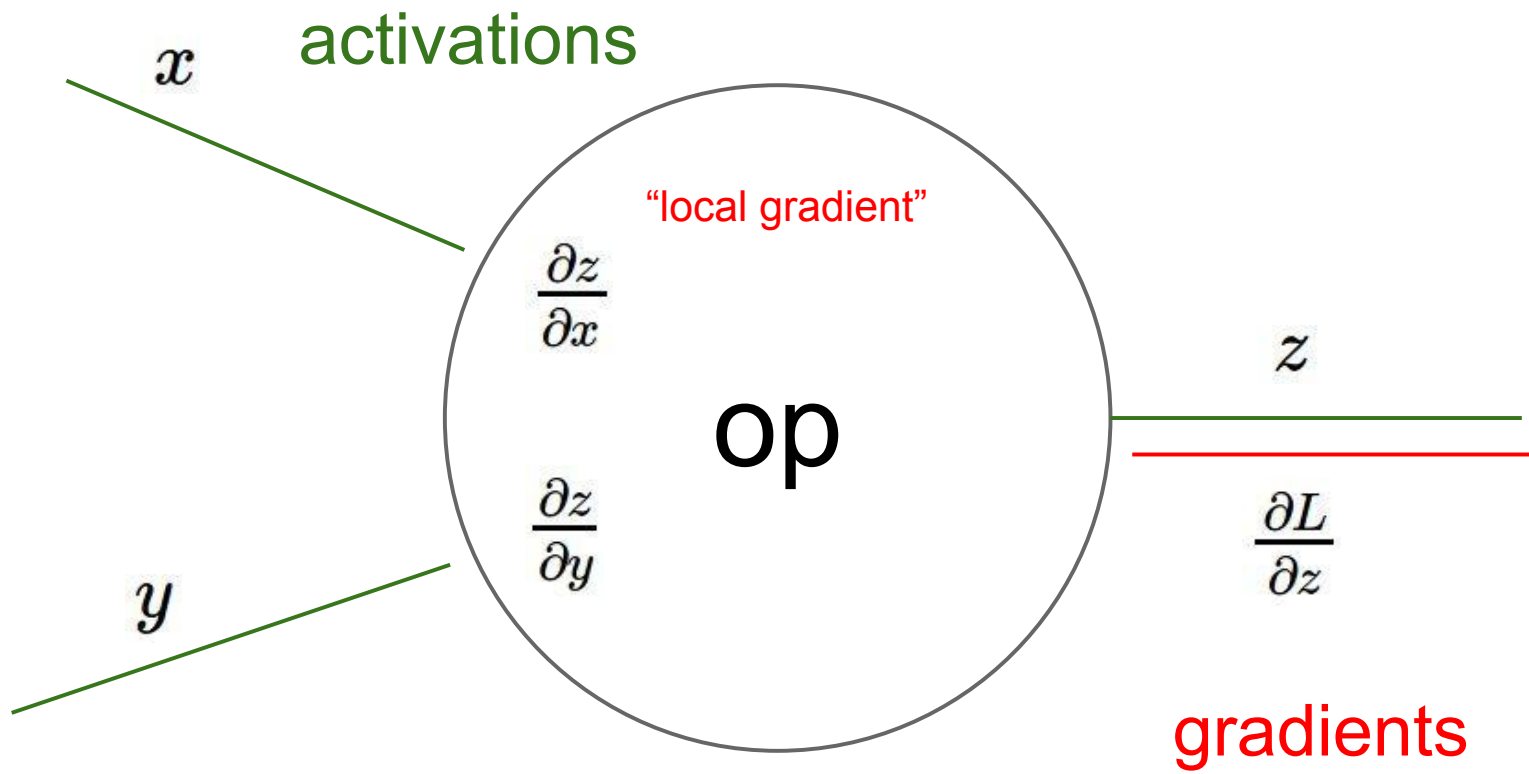https://colab.research.google.com/drive/1Il_bdC9of-_NQkW1jb1wcCpSheQr7j0T

https://github.com/ieee8023/NeuralNetwork-Examples/blob/master/mxnet/mxnet-vary-inputs-slideexamples.ipynb
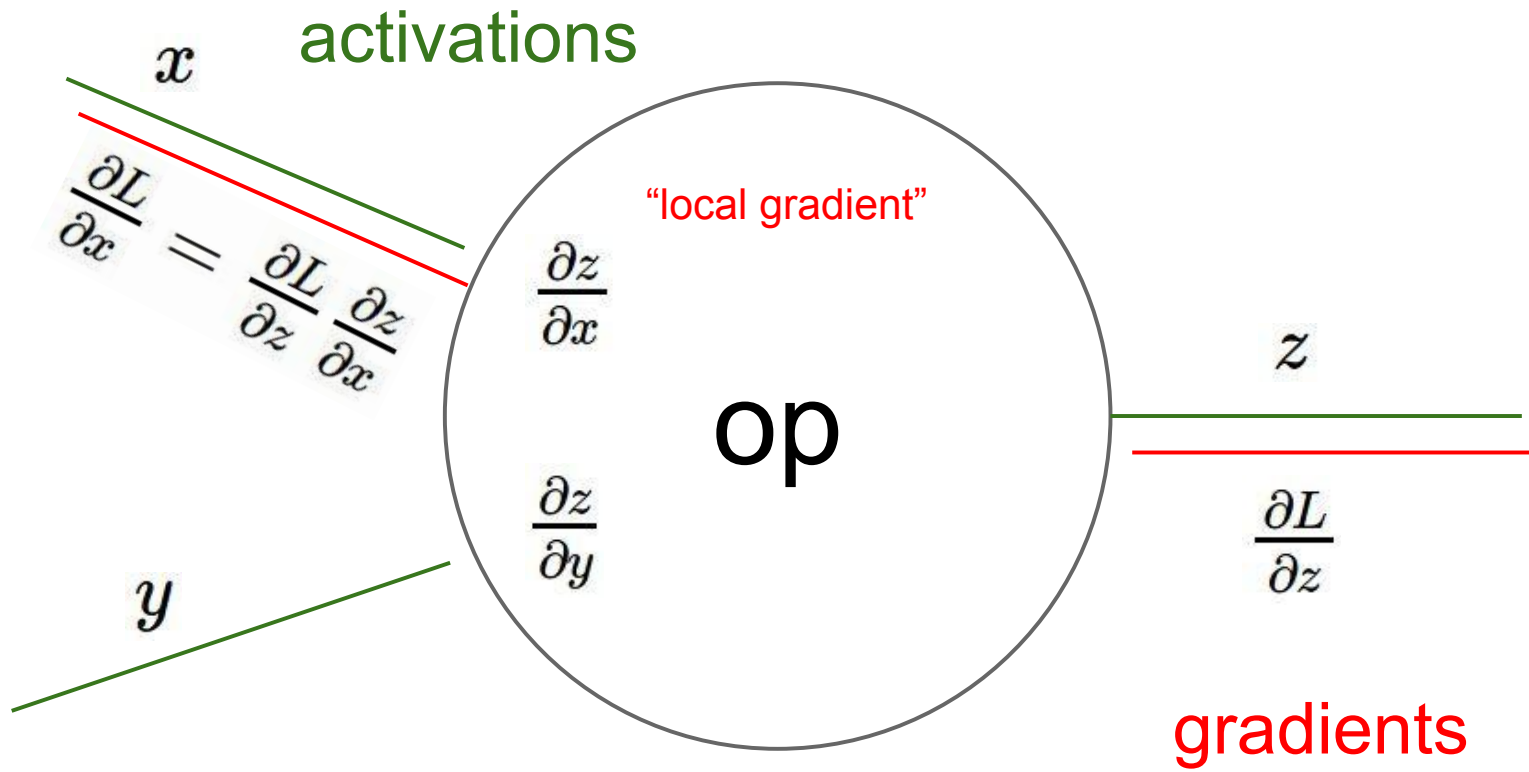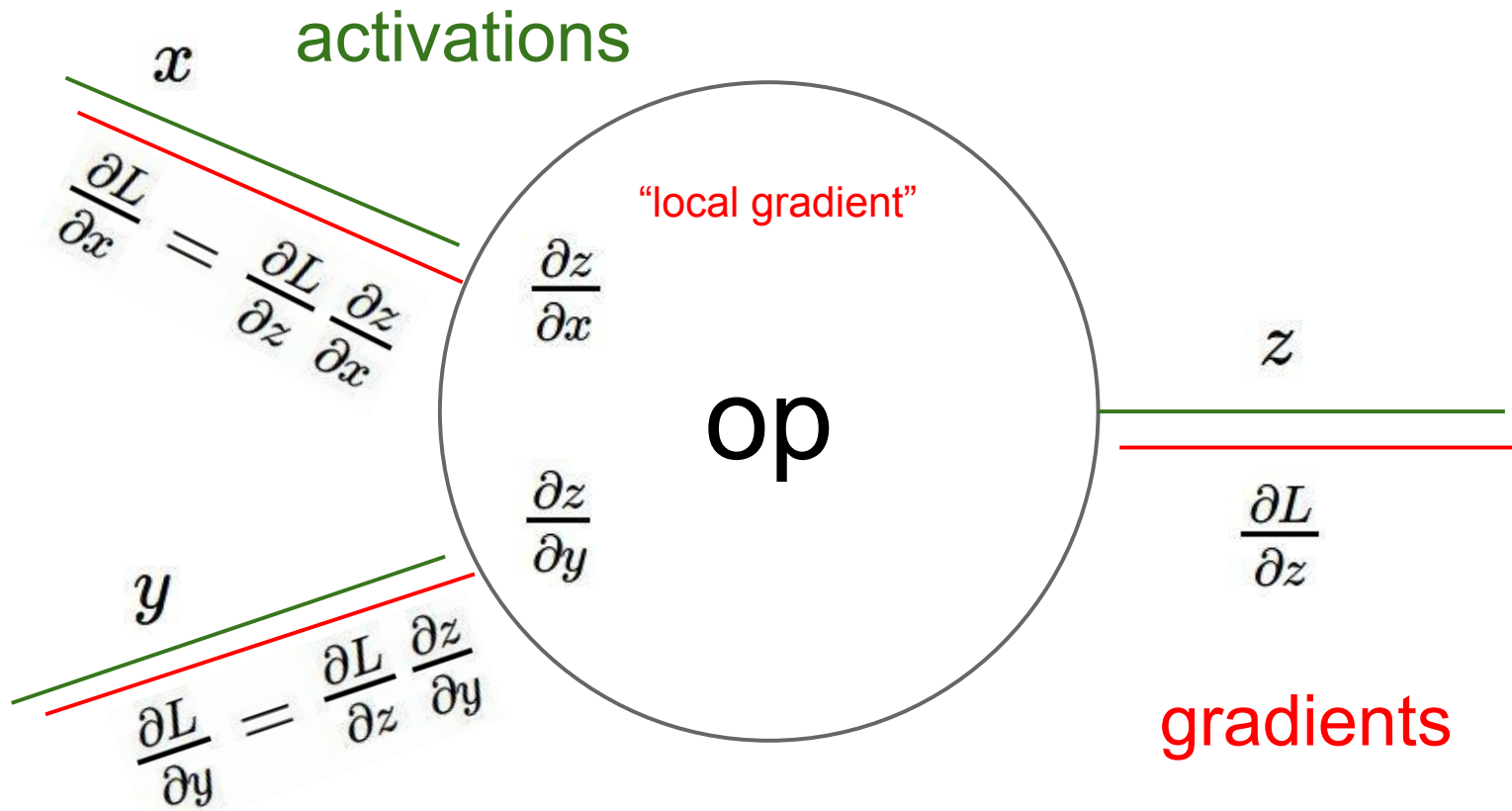
Distribute influence in the output of f to x and y

Distribute influence in the output of f to x and y

Distribute influence in the output of f to x and y

Distribute influence in the output of f to x and y

Distribute influence in the output of f to x and y

activations

$x$

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

"local gradient"

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

op

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

$z$

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

$y$

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

gradients

Distribute influence in the output of f to x and y

# Another example: Logistic (sigmoid) function

$$f\left(\sum_i w_i x_i + b\right) \qquad f(x) = \frac{1}{1 + e^{-x}}$$
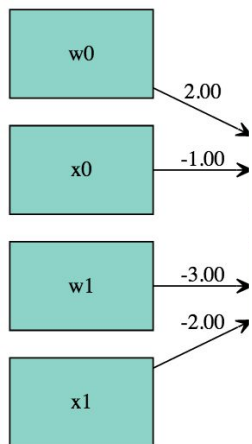
$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$\frac{1}{1 + e^{-\boxed{(w_0 x_0 + w_1 x_1 + w_2)}}}$$

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$\frac{1}{1 + \boxed{e^{-(w_0 x_0 + w_1 x_1 + w_2)}}}$$

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

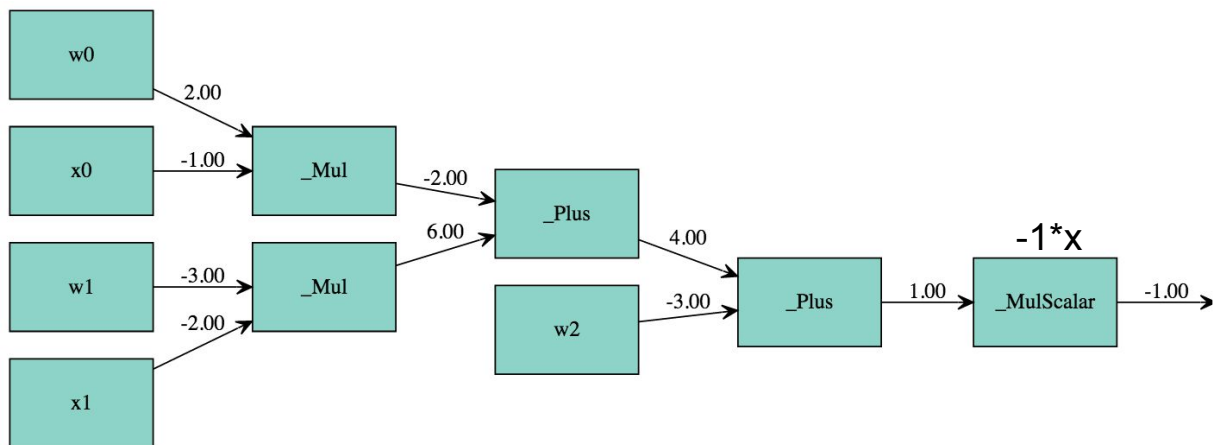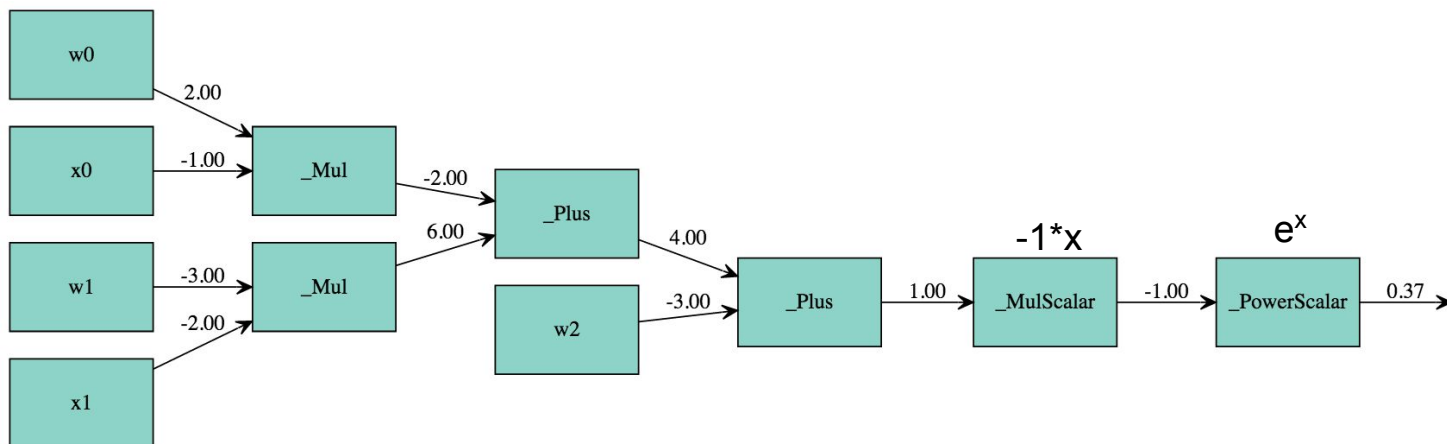$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$\frac{\partial f}{\partial f} = 1$$



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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Chain rule:

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

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

| | | |
|---|---|---|
| w0 | 2.00 | |
| x0 | -1.00 | _Mul → -2.00 |
| w1 | -3.00 | _Mul |
| x1 | -2.00 | |

_Plus → 4.00

w2 → -3.00

_Plus → 1.00 → _MulScalar → -1.00 → _PowerScalar → 0.37 → _PlusScalar → 1.37 → 1/x _DivScalar → 0.73

-.53    1.00

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

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

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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

1+x

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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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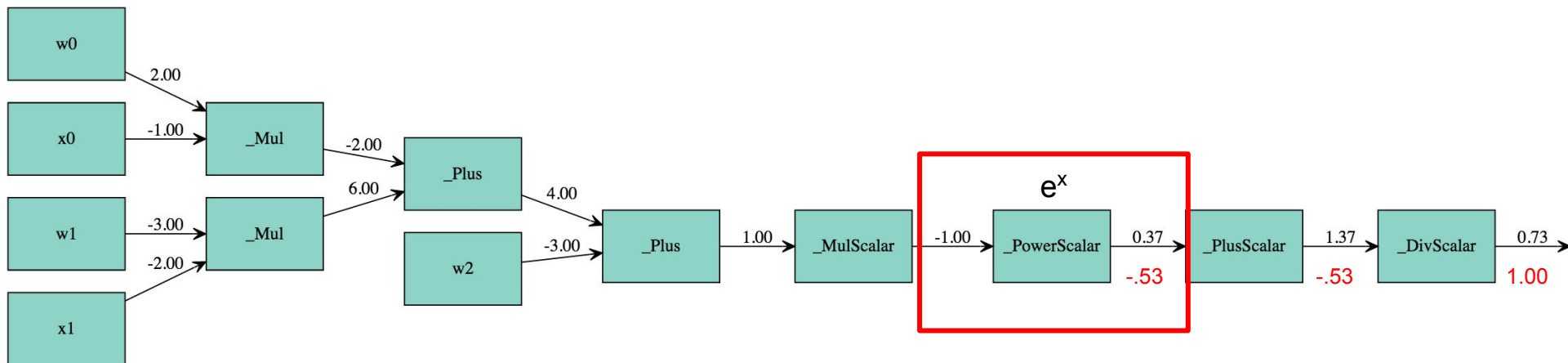
$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

w0

2.00

x0

-1.00

_Mul

-2.00

w1

-3.00

_Mul

6.00

_Plus

4.00

w2

-2.00

-3.00

_Plus

1.00

_MulScalar

-1.00

$e^x$

_PowerScalar

0.37

-.53

_PlusScalar

1.37

-.53

_DivScalar

0.73

1.00

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

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

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

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

$$\frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

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



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

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

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

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

$$\frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

(-1) * (-0.20) = 0.20

-1*x

| | | | | |
|---|---|---|---|---|
| w0 | | | | |
| x0 → _Mul | -2.00 | | | |
| | → _Plus | 4.00 | | |
| w1 → _Mul | 6.00 | | | |
| x1 | | | | |
| w2 → _Plus | 1.00 .20 → _MulScalar | -1.00 -.20 → _PowerScalar | 0.37 -.53 → _PlusScalar | 1.37 -.53 → _DivScalar | 0.73 1.00 |

w0 → 2.00
x0 → -1.00
w1 → -3.00
w1 → -2.00
w2 → -3.00

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

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

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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[local gradient] x [its gradient]
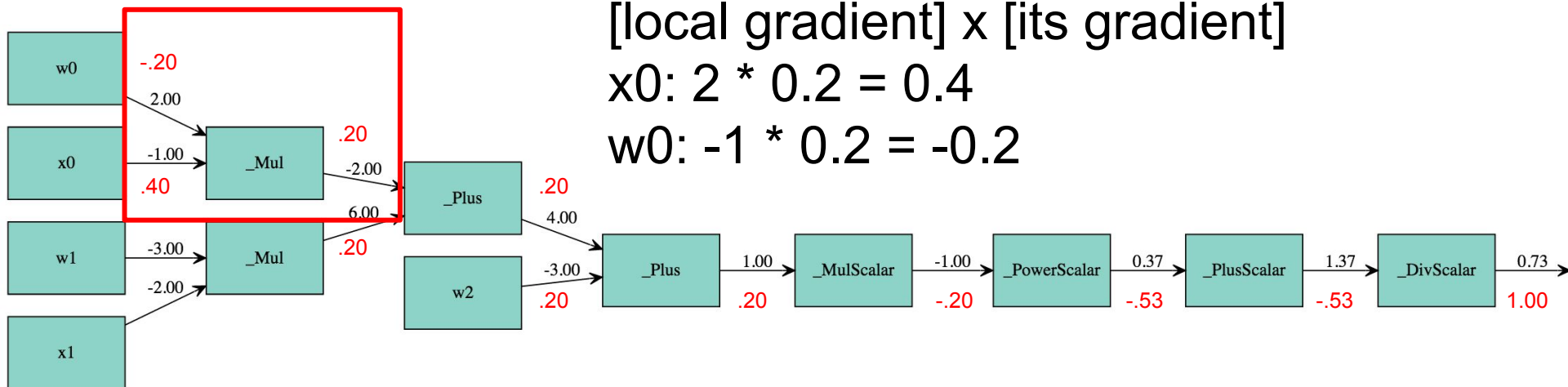1 * 0.2 = 0.2
1 * 0.2 = 0.2

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

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

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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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

$$\frac{1}{1+e^{-(w_0x_0+w_1x_1+w_2)}}$$

[local gradient] x [its gradient]
x0: 2 * 0.2 = 0.4
w0: -1 * 0.2 = -0.2

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

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

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

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

$$\frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



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

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

$$\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**add**: gradient distributor
**mul**: gradient scalar

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

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \quad \Big| \quad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$
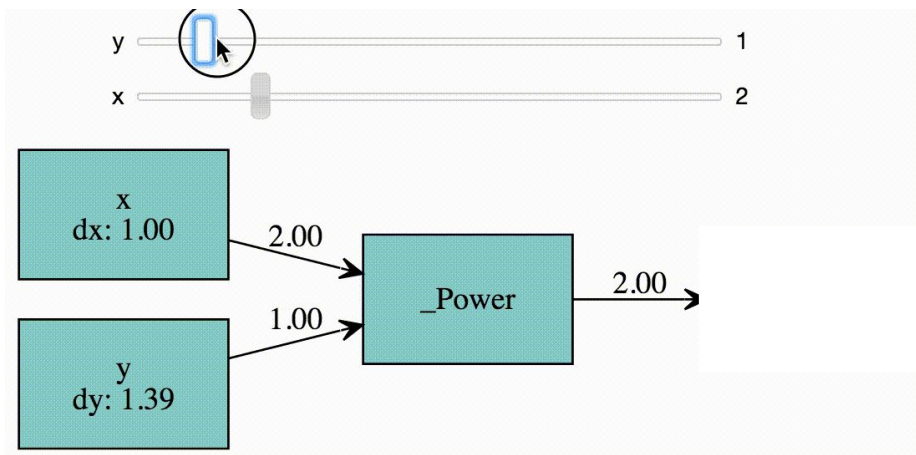
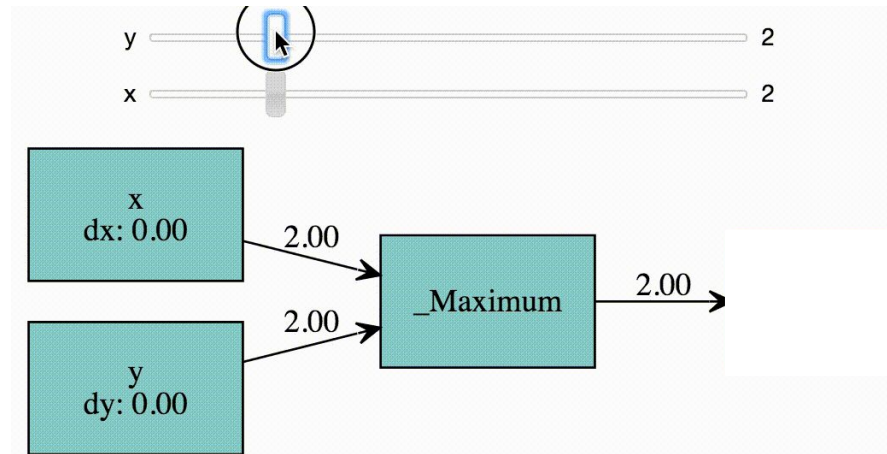**add**: gradient distributor
**mul**: gradient scalar

$x + y$

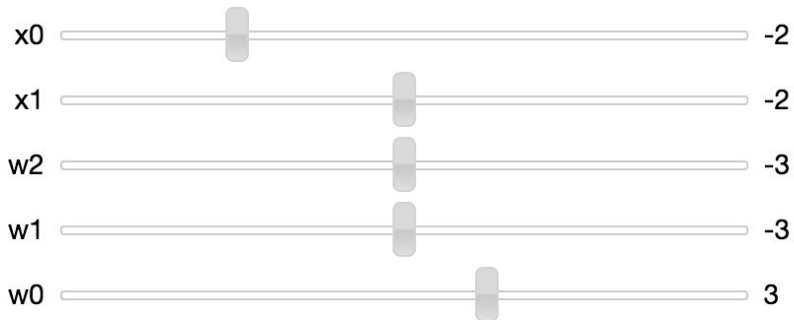$x \cdot y$
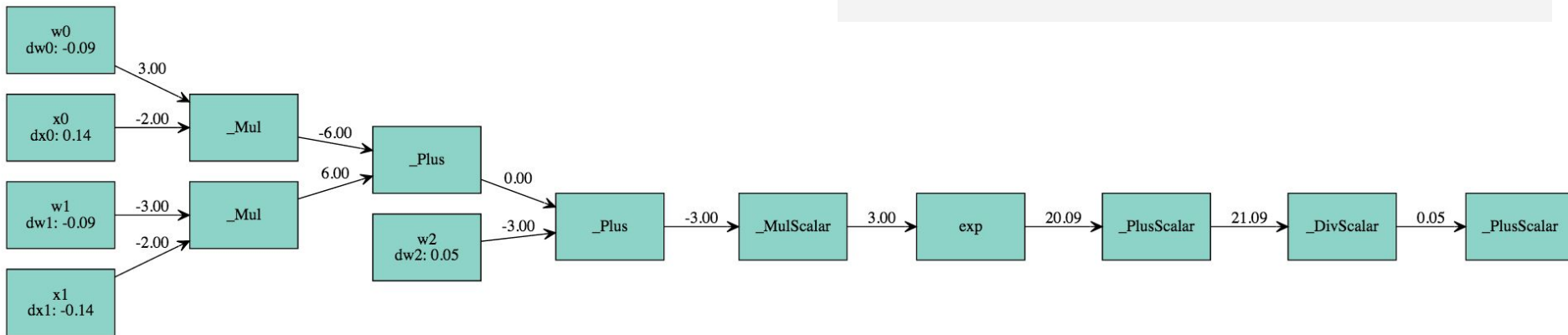
exp: gradient superscalar
max: gradient router

$x^y$

$max(x, y)$

# Interactive demo



1. Download IPython notebook from github
2. Install jupyter "pip install jupyter"
3. Run "jupyter notebook"
4. Navigate to mxnet-vary-inputs-slideexamples.ipynb
5. In the "Cell" menu select "Run All"

https://colab.research.google.com/drive/1Il_bdC9of-_NQkW1jb1wcCpSheQr7j0T

https://github.com/ieee8023/NeuralNetwork-Examples/blob/master/mxnet/mxnet-vary-inputs-slideexamples.ipynb