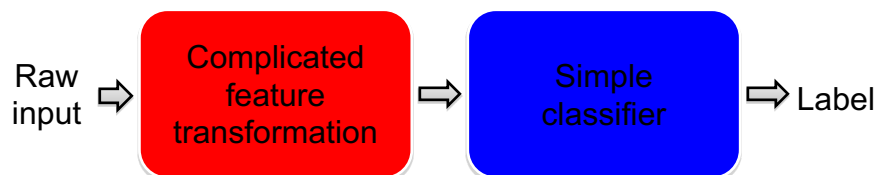


Multi-layer neural networks and backpropagation

Slides courtesy L. Lazebnik (Univ. of Illinois CS498) and others

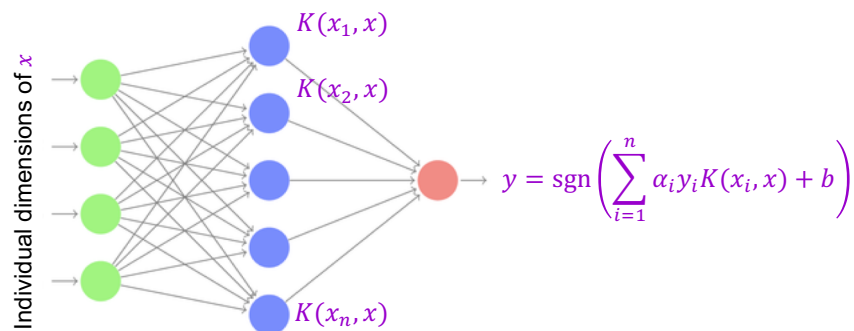
Beyond linear predictors

- To achieve good accuracy on challenging problems, we need to be able to train *nonlinear* models
- Traditional “shallow” approach:



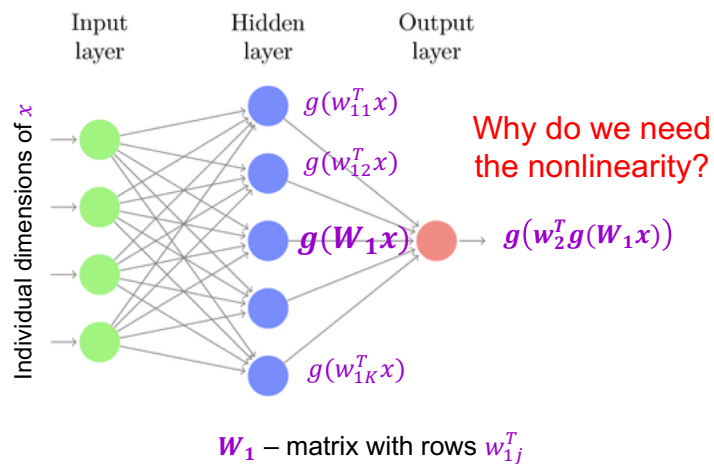
“Shallow” pipeline: Nonlinear SVM

- Perform a nonlinear mapping induced by kernel function, apply linear classifier
- Equivalently, compute kernel function value of input with every support vector, apply linear classifier



Two-layer neural network

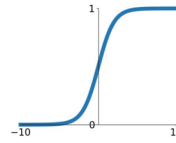
- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by a *nonlinearity*



Common nonlinearities

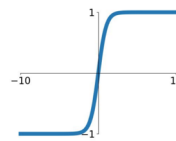
Sigmoid

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



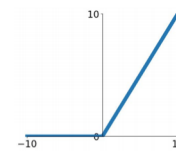
tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$



Source: [Stanford 231n](#)

Two-layer neural network

- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by a *nonlinearity*
- This gives a [universal function approximator](#)
 - But the hidden layer may need to be huge

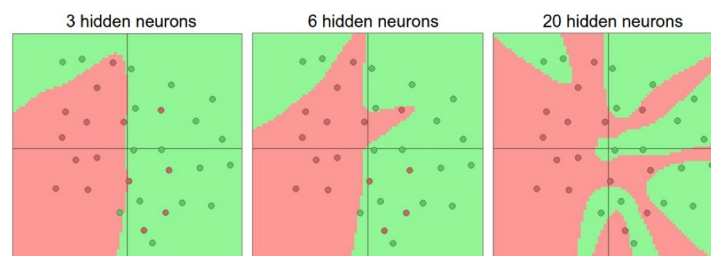
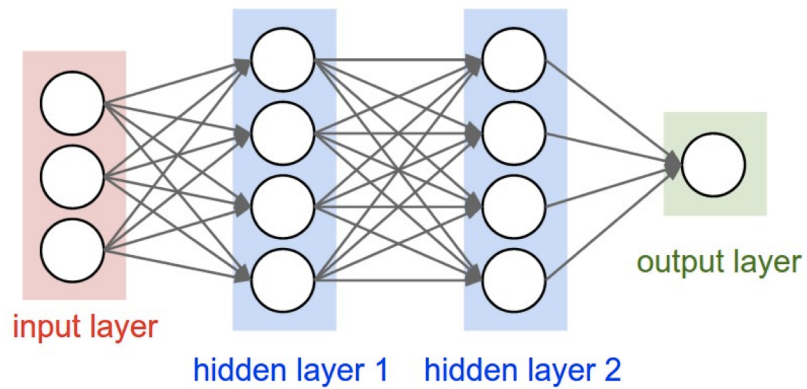


Figure source

Beyond two layers

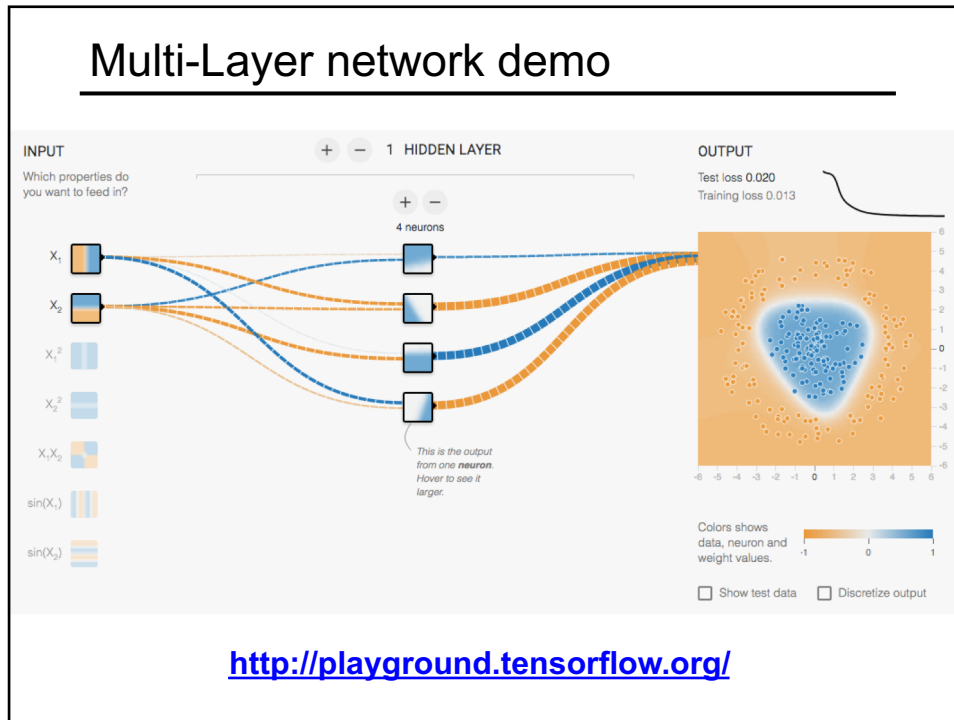


“Deep” pipeline

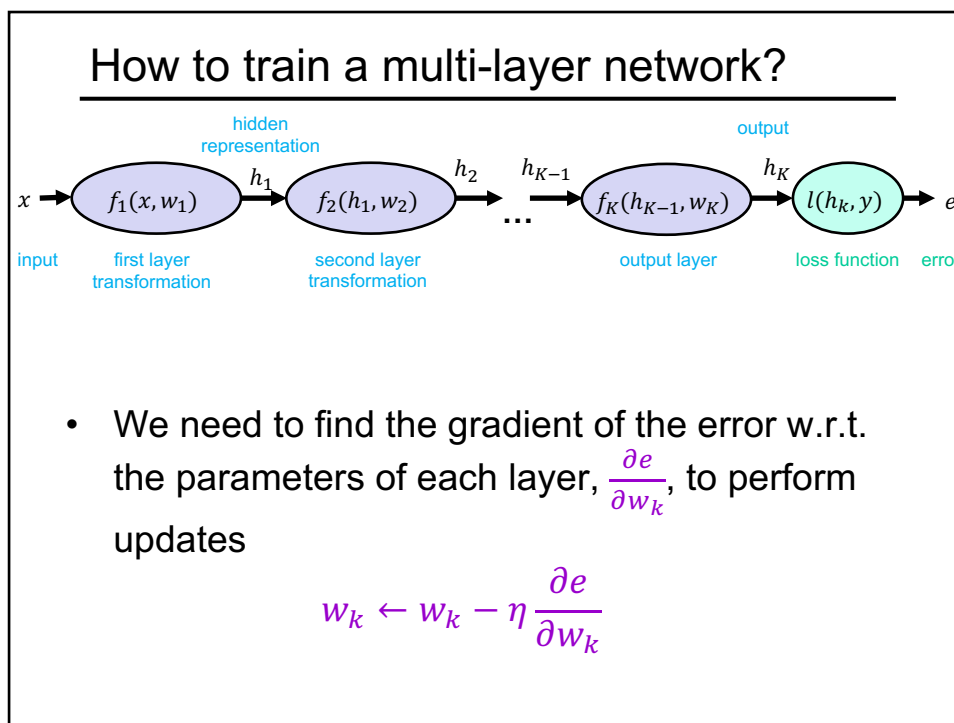


- Learn a *feature hierarchy*
- Each layer extracts features from the output of previous layer
- All layers are trained jointly

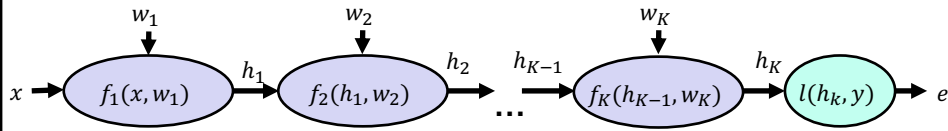
Multi-Layer network demo



How to train a multi-layer network?

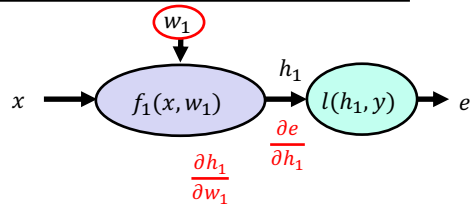


Computation graph



Chain rule

Let's start with $k = 1$



$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) =$$

Example: $e = (y - w_1^T x)^2$

$$h_1 = f_1(x, w_1) = w_1^T x$$

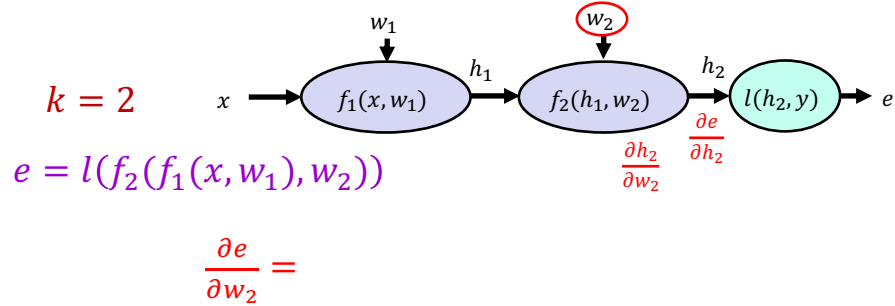
$$e = l(h_1, y) = (y - h_1)^2$$

$$\frac{\partial h_1}{\partial w_1} =$$

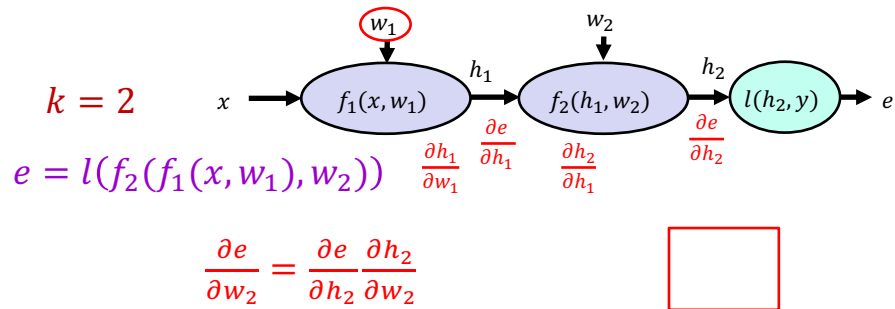
$$\frac{\partial e}{\partial h_1} =$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

Chain rule



Chain rule



Example: $e = -\log(\sigma(w_1^T x))$ (assume $y = 1$)

$$h_1 = f_1(x, w_1) = w_1^T x$$

$$h_2 = f_2(h_1) = \sigma(h_1)$$

$$e = l(h_2, 1) = -\log(h_2)$$

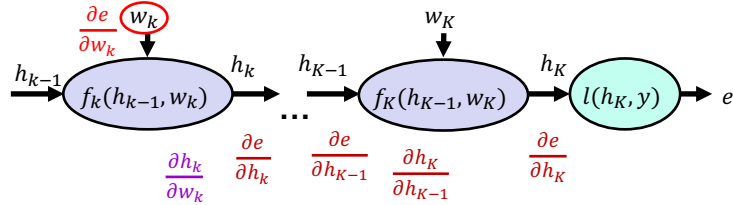
$$\frac{\partial h_1}{\partial w_1} =$$

$$\frac{\partial h_2}{\partial h_1} =$$

$$\frac{\partial e}{\partial h_2} =$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1} =$$

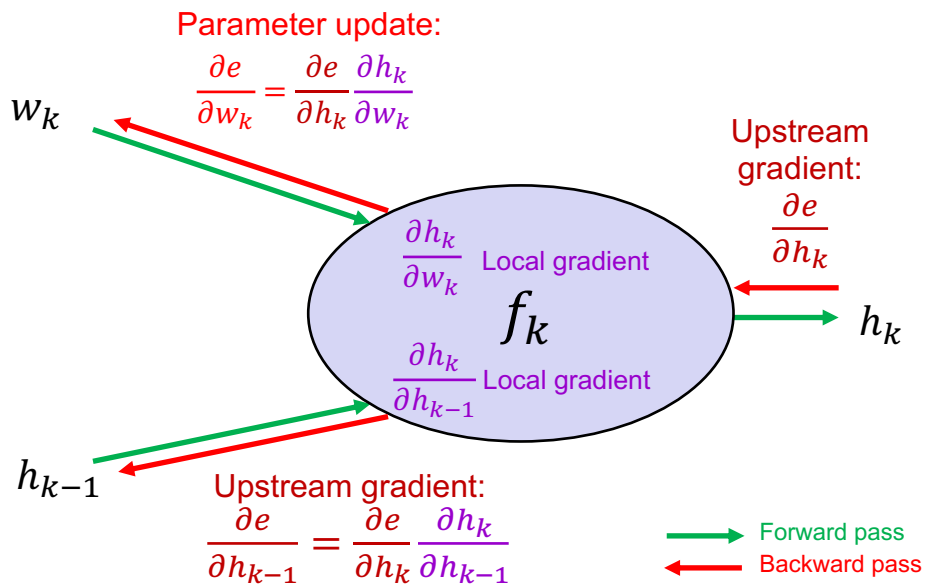
Chain rule



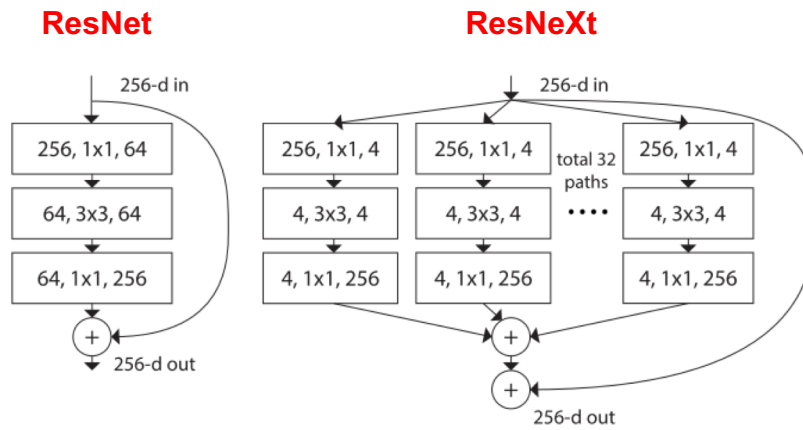
General case:

$$\frac{\partial e}{\partial w_k} = \underbrace{\frac{\partial e}{\partial h_K} \frac{\partial h_K}{\partial h_{k-1}} \cdots \frac{\partial h_{k+1}}{\partial h_k}}_{\text{Upstream gradient } \frac{\partial e}{\partial h_k}} \underbrace{\frac{\partial h_k}{\partial w_k}}_{\text{Local gradient}}$$

Backpropagation summary

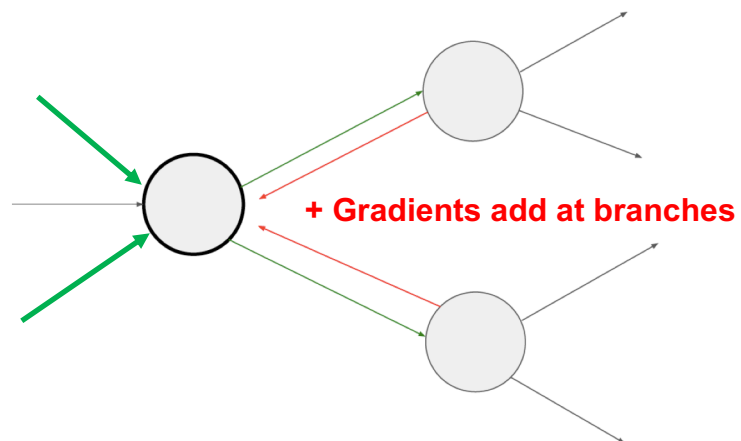


What about more general computation graphs?



[Figure source](#)

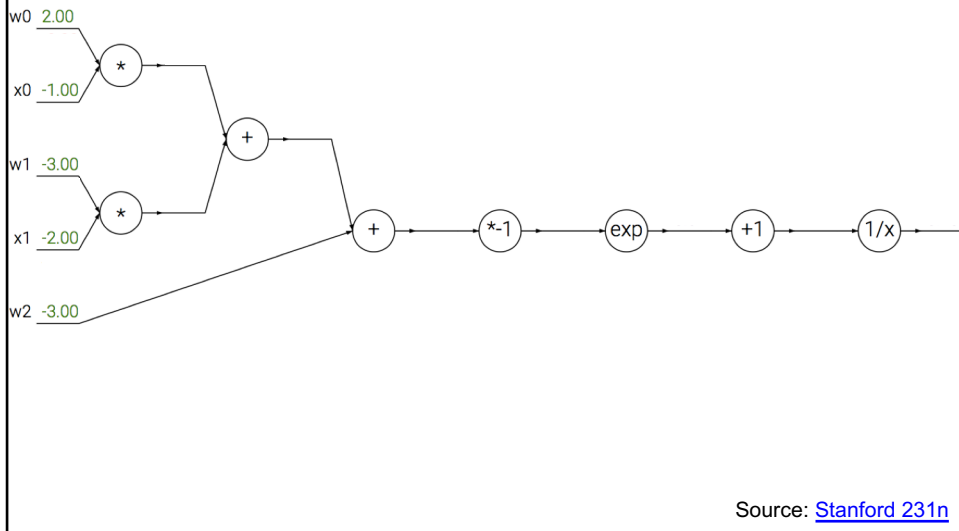
What about more general computation graphs?



Source: [Stanford 231n](#)

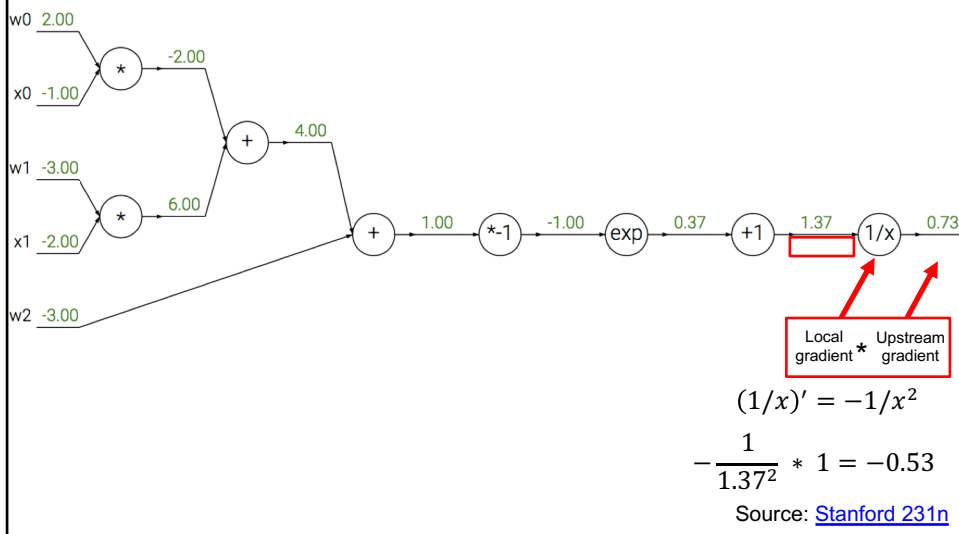
A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



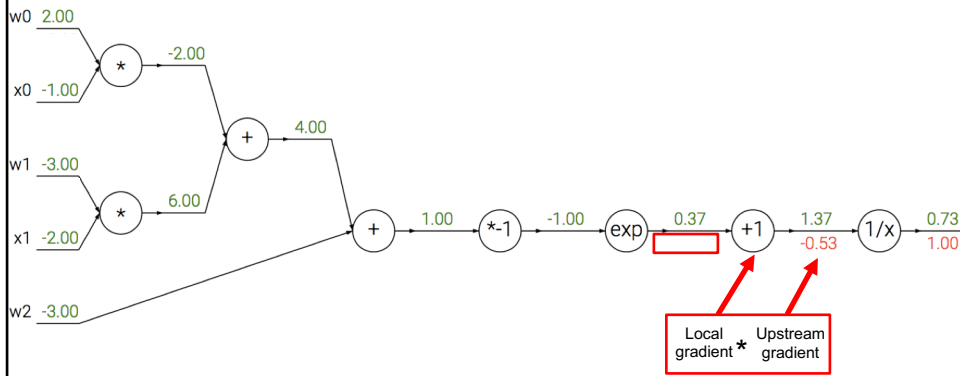
A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



A detailed example

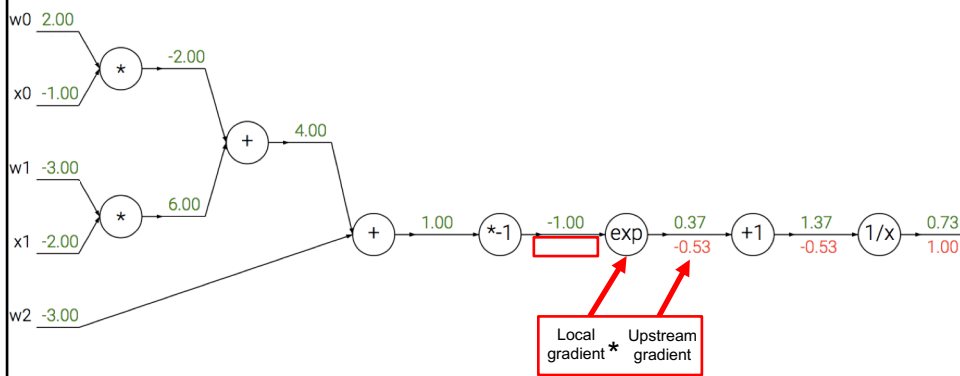
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$

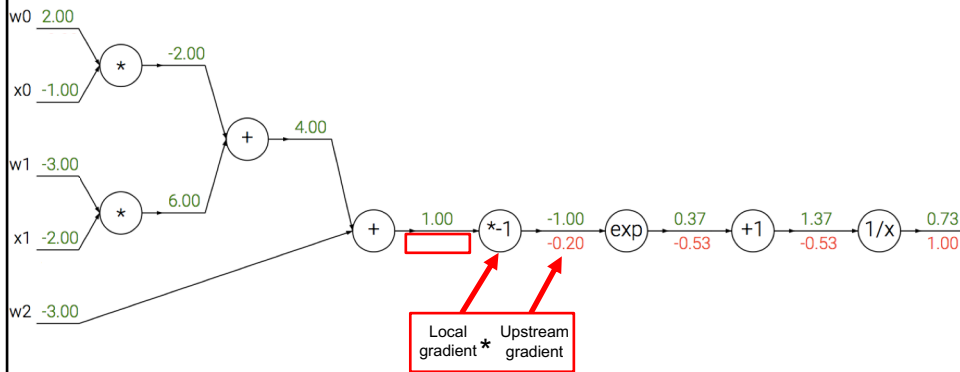


$$\exp(-1) * (-0.53) = -0.20$$

Source: [Stanford 231n](#)

A detailed example

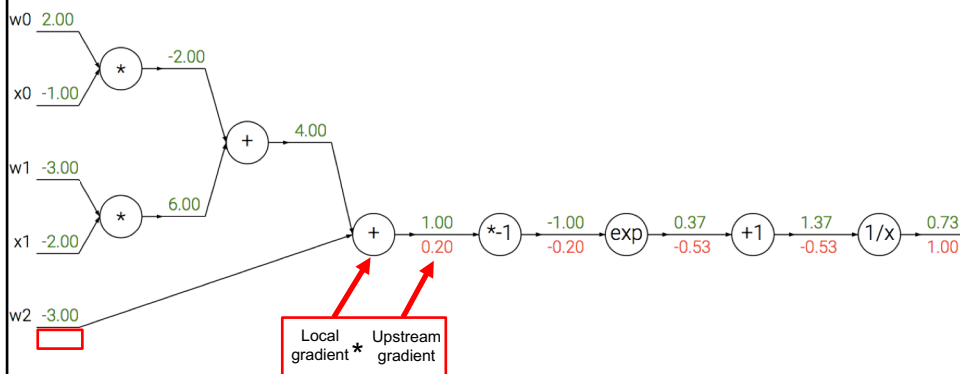
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

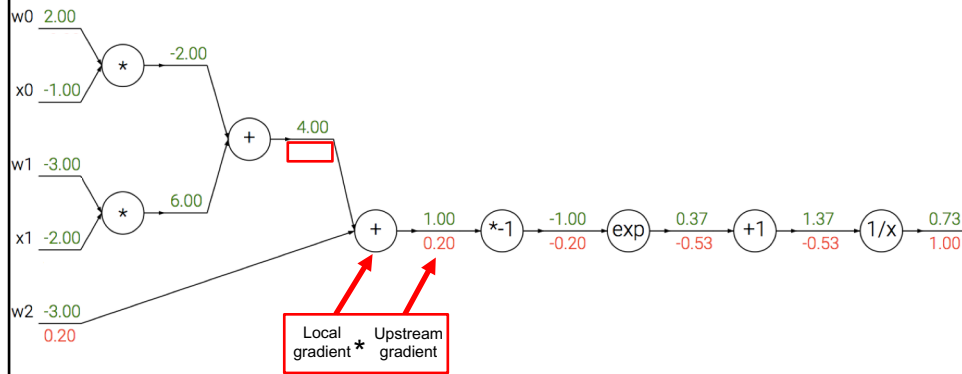
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

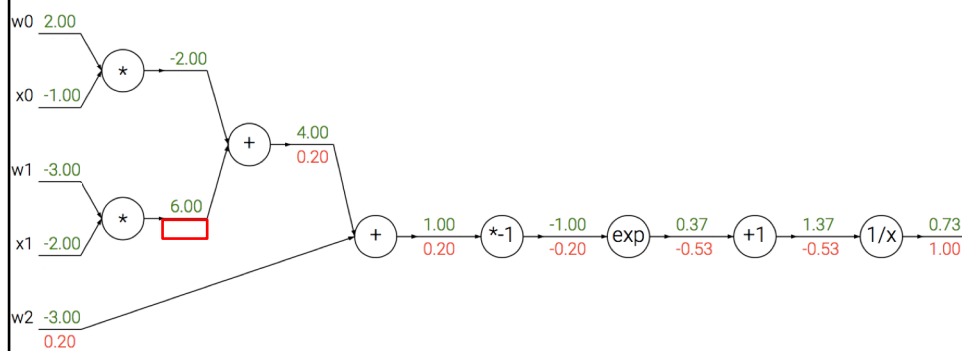
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

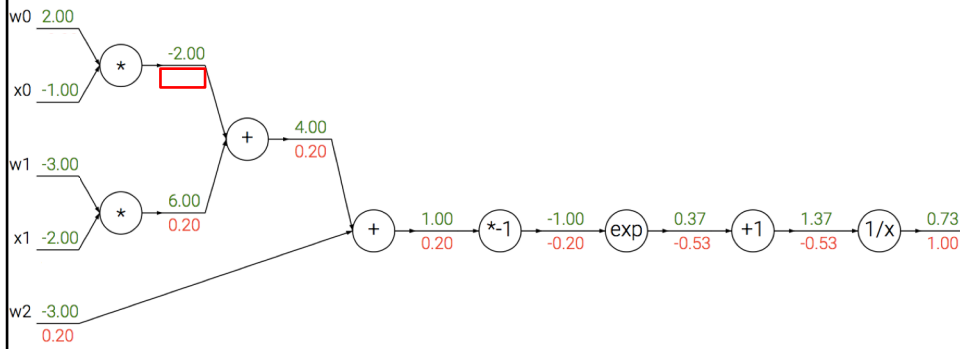
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

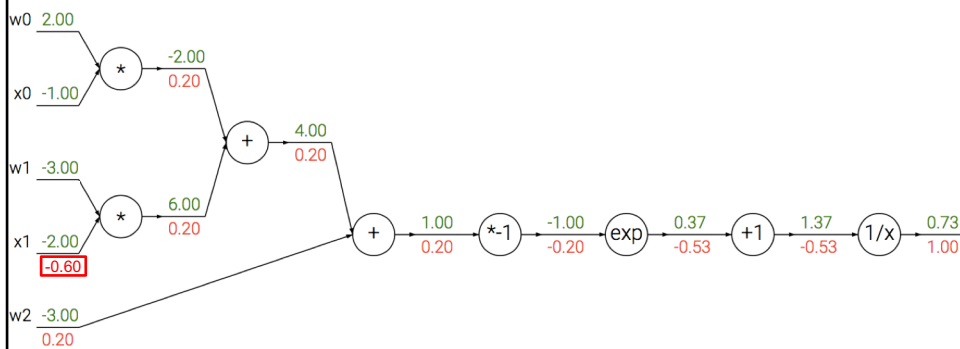
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

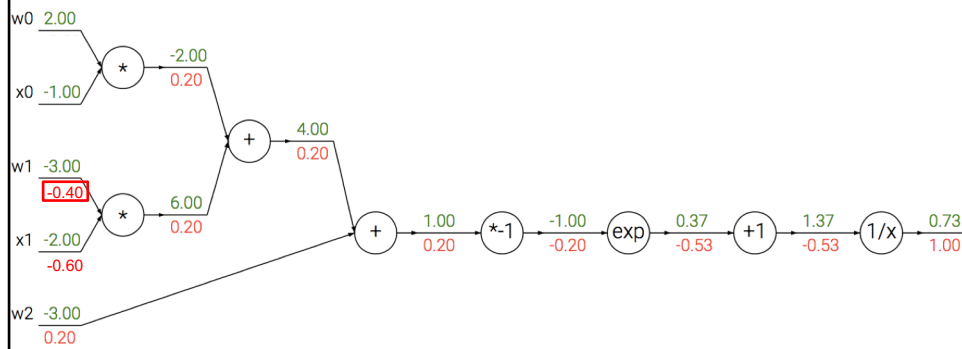
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

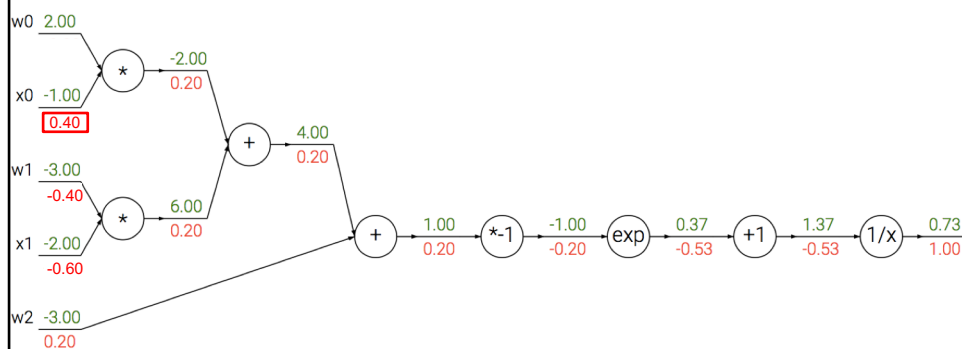
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

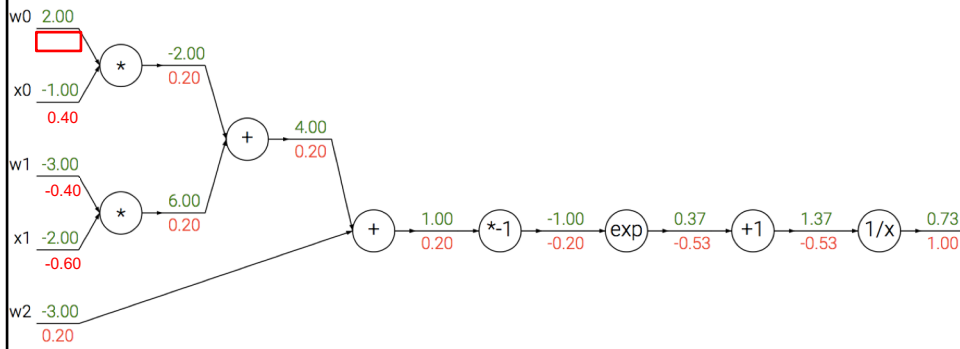
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

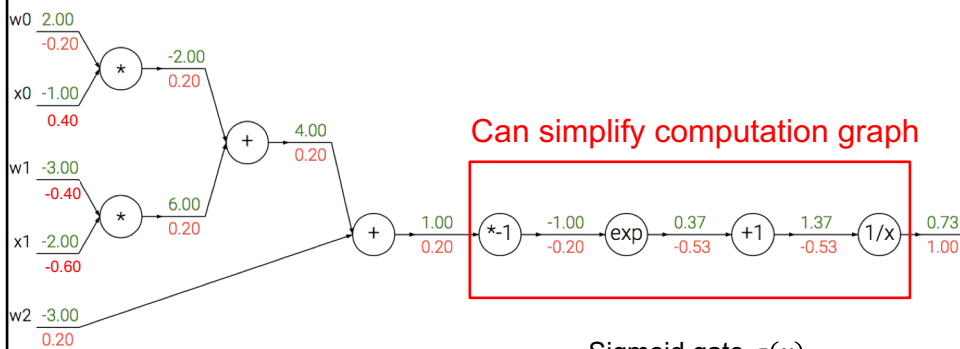
$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



Source: [Stanford 231n](#)

A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0x_0 + w_1x_1 + w_2)]}$$



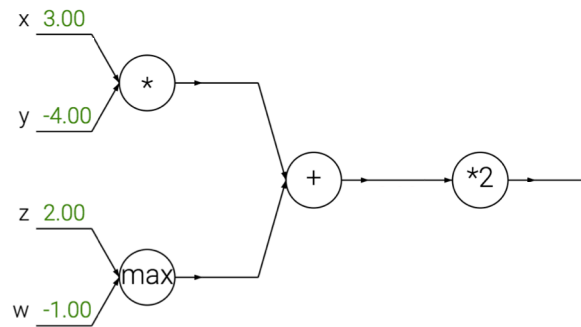
Sigmoid gate $\sigma(x)$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$\sigma(1)(1 - \sigma(1)) = 0.73 * (1 - 0.73) = 0.20$$

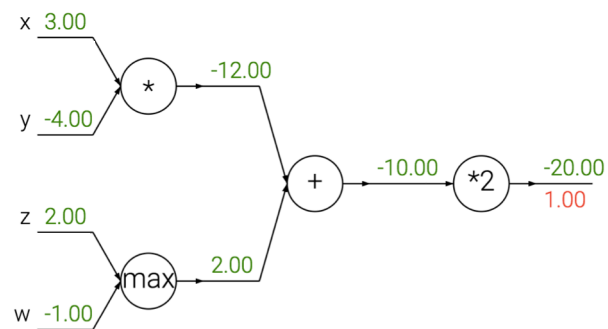
Source: [Stanford 231n](#)

Patterns in gradient flow



Source: [Stanford 231n](#)

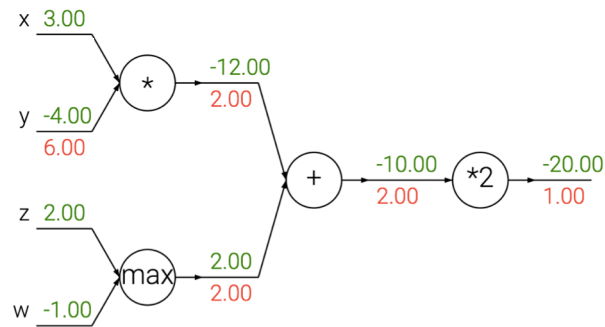
Patterns in gradient flow



Add gate: "gradient distributor"

Source: [Stanford 231n](#)

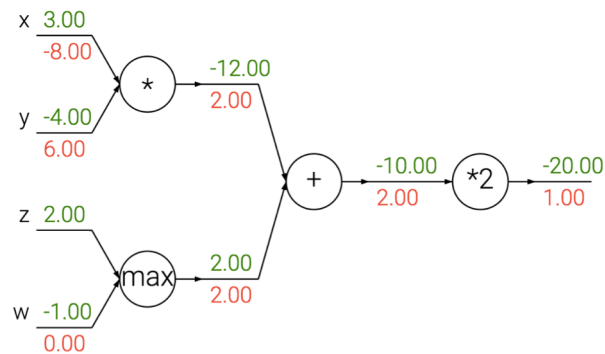
Patterns in gradient flow



Add gate: "gradient distributor"
Multiply gate: "gradient switcher"

Source: [Stanford 231n](#)

Patterns in gradient flow



Add gate: "gradient distributor"
Multiply gate: "gradient switcher"
Max gate: "gradient router"

Source: [Stanford 231n](#)

Dealing with vectors

$$\frac{\partial z}{\partial x} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_1}{\partial x_M} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_N}{\partial x_1} & \dots & \frac{\partial z_N}{\partial x_M} \end{pmatrix}$$

$N \times M$
Jacobian

x $\xrightarrow{\text{green}}$ $f(x)$ $\xrightarrow{\text{green}}$ z

$1 \times M$ $\frac{\partial e}{\partial x} = \frac{\partial e}{\partial z} \frac{\partial z}{\partial x}$ $1 \times N$

$1 \times M$ $1 \times N$ $N \times M$ $1 \times N$

Simple case: Elementwise operation

$$\frac{\partial z}{\partial x} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_1}{\partial x_M} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_M}{\partial x_1} & \dots & \frac{\partial z_M}{\partial x_M} \end{pmatrix}$$

$M \times M$
Jacobian

x $\xrightarrow{\text{green}}$ $f(x) = \max(0, x)$ $\xrightarrow{\text{green}}$ z

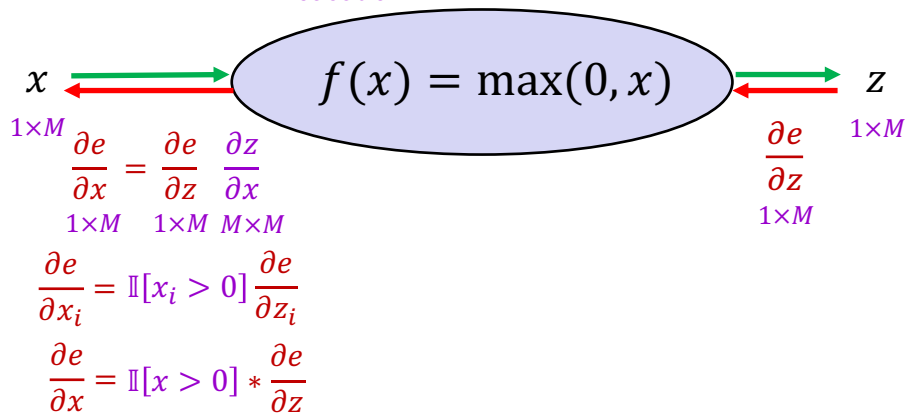
$1 \times M$ $\frac{\partial e}{\partial z}$ $1 \times M$

$1 \times M$

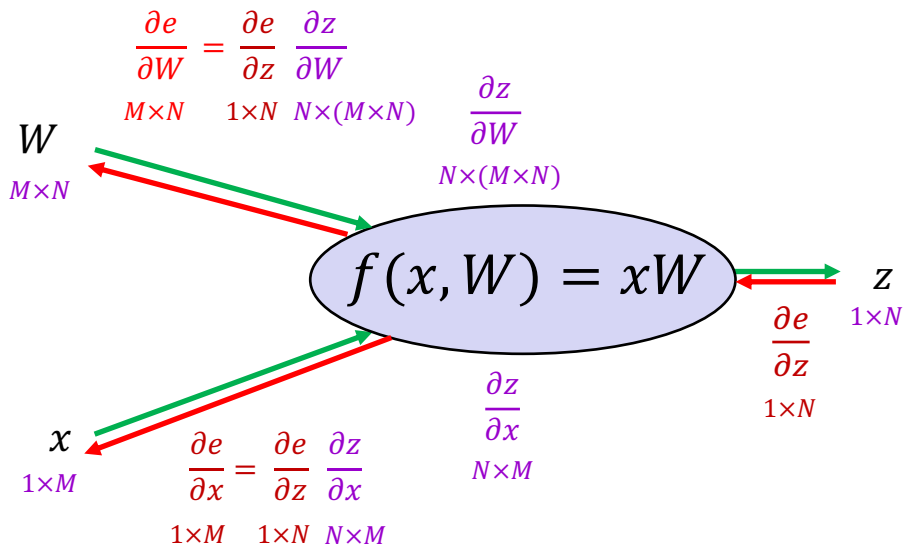
Simple case: Elementwise operation

$$\frac{\partial z}{\partial x} = \begin{pmatrix} \mathbb{I}[x_1 > 0] & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \mathbb{I}[x_M > 0] \end{pmatrix}$$

$M \times M$
Jacobian



Matrix-vector multiplication



General tips

- Derive error signal (upstream gradient) directly, avoid explicit computation of huge local derivatives
- Write out expression for a single element of the Jacobian, then deduce the overall formula
- Keep consistent indexing conventions, order of operations
- Use dimension analysis
- **Useful resource:** see Lecture 4 of [Stanford 231n](#) and associated links in the syllabus