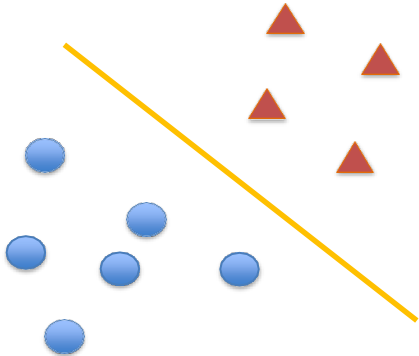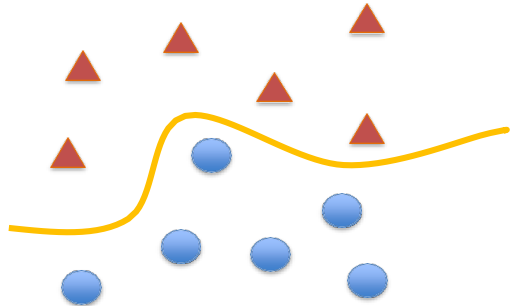# L8 – Neural networks

Linear Model
(Perceptron)

Need a non-linear model!

# Non-linear models: Neural networks

What should we do in case the target is not linearly separable?

- Use **multiple linear functions** and combine them using **non-linear activation functions**

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

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

  **activation function**

  $x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$

# Non-linear models: Neural networks

- A "neuron" computes a linear function $f = Wx$

- Neural network is a nesting of 'functions'
  - 2-layers: $f = W_2 \boxed{\max(0, W_1 x)}$ ← **activation function**
  - 3-layers: $f = W_3 \max(0, W_2 \max(0, W_1 x))$
  - 4-layers: $f = W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x)))$
  - 5-layers: $f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$
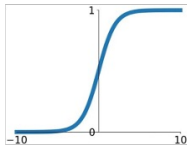  - … up to hundreds of layers

## Activation functions

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



**tanh**
$\tanh(x)$



**ReLU**
$\max(0, x)$

## 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"

Neural network diagram with input layer $x_1$, $x_2$, $x_3$; first hidden layer $f(W_{0,0}x + b_{0,0})$, $f(W_{0,1}x + b_{0,1})$, $f(W_{0,2}x + b_{0,2})$, $f(W_{0,3}x + b_{0,3})$; second hidden layer $f(W_{1,0}x + b_{1,0})$, $f(W_{1,1}x + b_{1,1})$, $f(W_{1,2}x + b_{1,2})$; output layer $f(W_{2,0}x + b_{2,0})$.

- Split your data



Other splits are also possible (e.g., 80% / 10% / 10%)

# Training recipe for MLP



Inputs → Neural Network → Outputs ↔ Targets

**Are these reasonably close?**

large distance!

target

bad prediction

prediction

small distance!

target

good prediction

- Select a suitable loss function

**Binary cross entropy loss**

$$\mathcal{L}_{\text{BCE}} = -\frac{1}{N} \sum_{i=1}^{N} [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

- N is the total number of samples
- $y_i$ is the true label for the i-th sample (either 0 or 1).
- $\hat{y}_i$ is the predicted probability for the i-th sample.



Yes! (0.8)

No! (0.2)

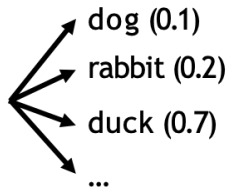The network predicts the probability of the input belonging to the "yes" class

**Cross entropy loss**

$$\mathcal{L}_{\mathrm{CE}} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} y_{i,c} \log(\hat{y}_{i,c})$$

- N is the total number of samples
- $y_{i,c}$ is a binary indicator (1 if the true class of sample i is c, 0 otherwise
- $\hat{y}_i$ is the predicted probability for the i-th sample.
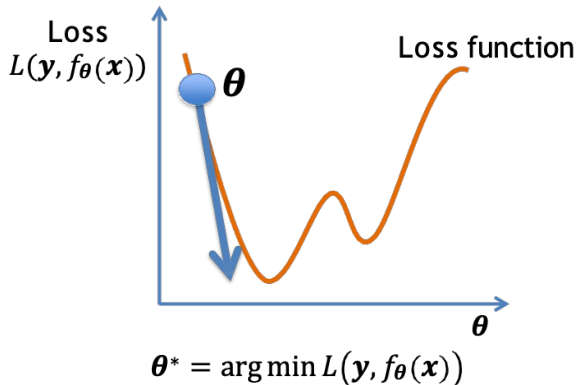- C is the number of classes



dog (0.1)
rabbit (0.2)
duck (0.7)
...

# Training a neural net

- Minimize: $L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$ w.r.t. $\boldsymbol{\theta}$

- We use gradient-based optimization: **gradient descent**



Loss
$L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

Loss function

$\boldsymbol{\theta}$

$\boldsymbol{\theta}^* = \arg\min L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

Iterative update:

Learning rate

$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

# Training a neural net

- Given inputs $x$ and targets $y$

- Given multi-layer NN

  - Need to propagate gradients from end to first layer ($W_1$)

- Backpropagation: Use chain rule to compute gradients

The flow of computations in a neural network goes in two ways:

**1. Left-to-right**: This is referred to as *forward propagation,* which results in computing the output of the network

**2. Right-to-left**: This is referred to as *back propagation*, which results in computing the gradients (or derivatives) of the parameters in the network
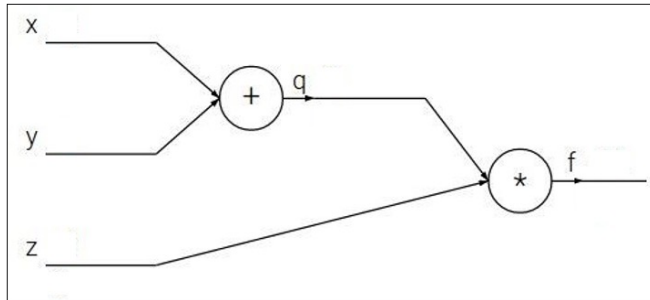
Backpropagation: a simple example

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

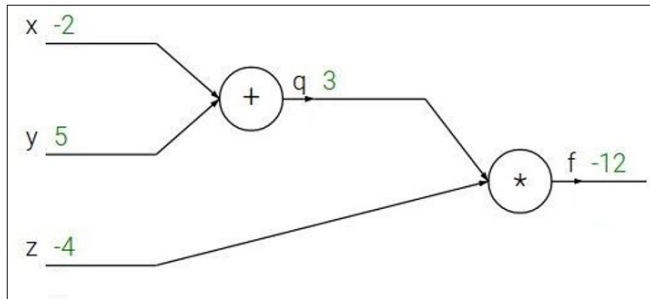Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$
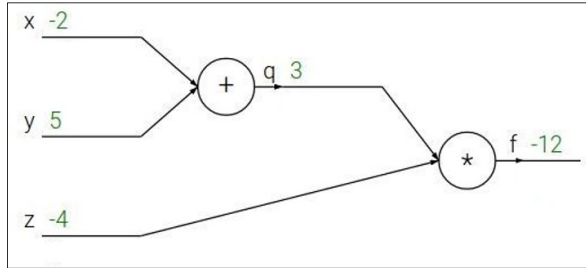
e.g. x = -2, y = 5, z = -4

Backpropagation: a simple example

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

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

# Backpropagation

Backpropagation: a simple example

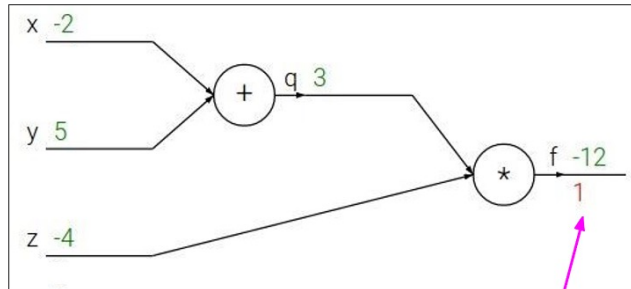$$f(x,y,z) = (x+y)z$$

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



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

Backpropagation: a simple example

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

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



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

# Backpropagation

Backpropagation: a simple example
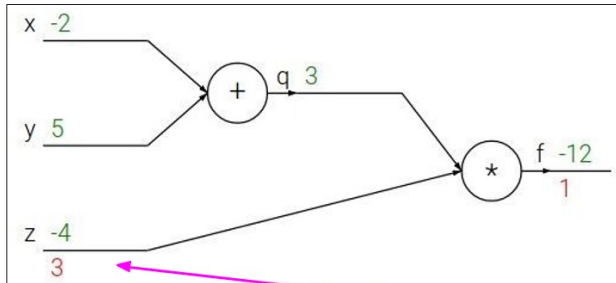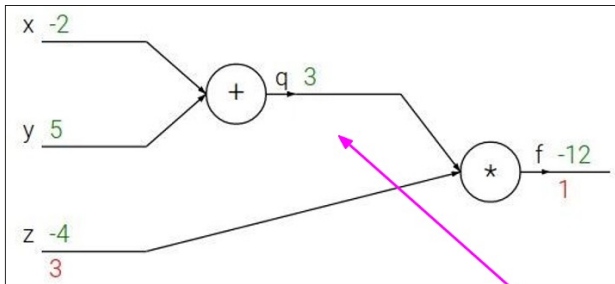
$$f(x, y, z) = (x + y)z$$

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



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

# Backpropagation
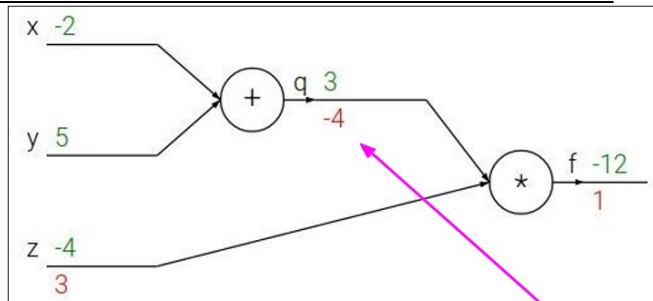
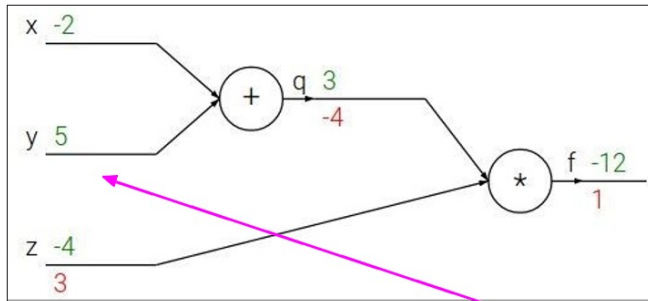Backpropagation: a simple example

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

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



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

Backpropagation: a simple example

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

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:

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

Backpropagation: a simple example

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

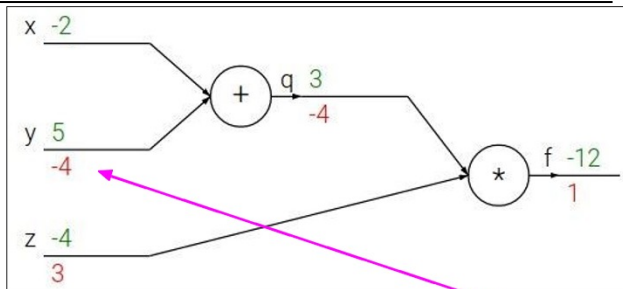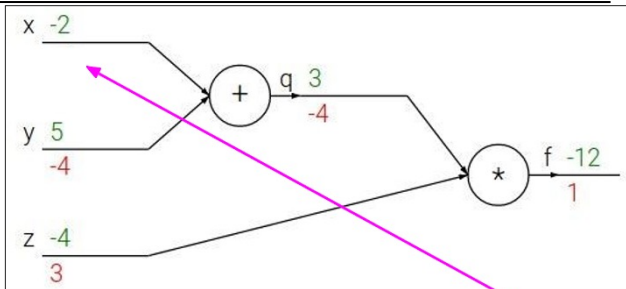e.g. x = -2, y = 5, z = -4

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

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

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



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

Chain rule:

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

Upstream gradient  Local gradient

# Backpropagation

Backpropagation: a simple example

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

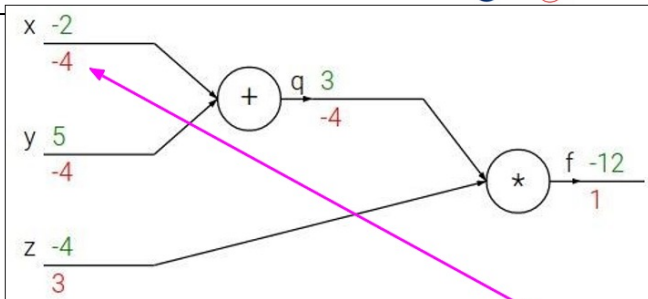e.g. x = -2, y = 5, z = -4

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

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

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



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

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$
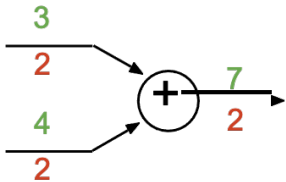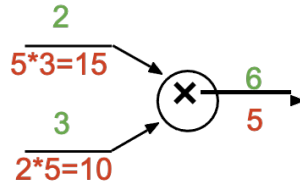
Upstream gradient    Local gradient

# Backpropagation

Backpropagation: a simple example

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

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:

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

Upstream gradient    Local gradient

## Patterns in gradient flow

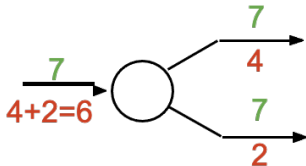**add** gate: gradient distributor



3
2
4
2
+
7
2
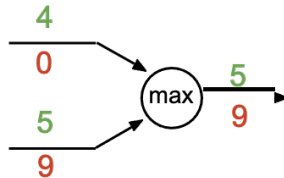
**mul** gate: "swap multiplier"



2
5*3=15
3
2*5=10
×
6
5

**copy** gate: gradient adder



7
4+2=6
7
4
7
2

**max** gate: gradient router



4
0
5
9
max
5
9

If we have a linear combination

$$W \cdot x = q$$

W matrix
x matrix or vector

Using the chain rule we can compute the gradients:

$$\nabla f_W = \nabla f_q \cdot x^T$$

$$\nabla f_x = W^T \cdot \nabla f_q$$

At test time / validation time
- Perform the forward pass and get the probabilities

- Binary classification: Threshold the probability to find the predicted class

$$\begin{cases} \text{if } p(y = 1|x; \theta) < 0.5 \text{ predict class } 0 \\ \text{if } p(y = 1|x; \theta) \geq 0.5 \text{ predict class } 1 \end{cases}$$
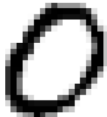
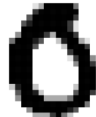Classify digits on the MNIST dataset
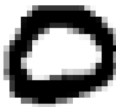
**TRAIN set**

[0]  [0]  [0]  [0]  [0]



**TEST set**

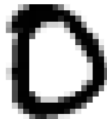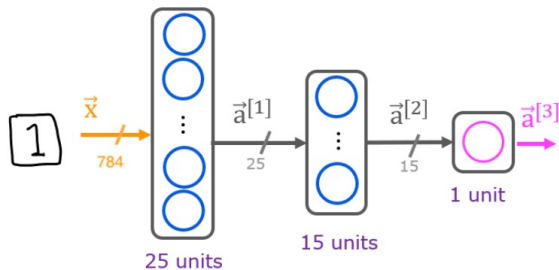[1]  [0]  [0]  [0]  [0]

Implement this architecture and compute using
backpropagation the gradients of the parameters



- Train the network on the training set
- Compute the accuracy on the test set