# Generic form of Feedforward Neural Networks

# Designing feedforward neural networks

$x_1 \rightarrow h_1^{(0)}$

$x_2 \rightarrow h_2^{(0)}$

$a_1^{(1)}$ $g^{(1)}$ $h_1^{(1)}$

$a_2^{(1)}$ $g^{(1)}$ $h_2^{(1)}$

$a_3^{(1)}$ $g^{(1)}$ $h_3^{(1)}$

First layer: $\mathbb{R}^2 \rightarrow \mathbb{R}^3$
$\boldsymbol{h}^{(1)} = g^{(1)}\big(\boldsymbol{a}^{(1)}\big)$
$\boldsymbol{a}^{(1)} = W^{(1)}\boldsymbol{h}^{(0)}$
$W^{(1)} \in \mathbb{R}^{3\times2}, \boldsymbol{a}^{(1)}, \boldsymbol{h}^{(1)} \in \mathbb{R}^3$

$a_1^{(2)}$ $g^{(2)}$ $h_1^{(2)}$

$a_2^{(2)}$ $g^{(2)}$ $h_2^{(2)}$

Second layer: $\mathbb{R}^3 \rightarrow \mathbb{R}^2$
$\boldsymbol{h}^{(2)} = g^{(2)}\big(\boldsymbol{a}^{(2)}\big)$
$\boldsymbol{a}^{(2)} = W^{(2)}\boldsymbol{h}^{(1)}$
$W^{(2)} \in \mathbb{R}^{2\times3}, \boldsymbol{a}^{(2)}, \boldsymbol{h}^{(2)} \in \mathbb{R}^2$

$a_1^{(3)}$ $g^{(3)}$ $h_1^{(3)} \leftarrow y_1$

$a_1^{(3)}$ $g^{(3)}$ $h_2^{(3)} \leftarrow y_2$

Final layer: $\mathbb{R}^2 \rightarrow \mathbb{R}^2$
$\boldsymbol{h}^{(3)} = g^{(3)}\big(\boldsymbol{a}^{(3)}\big)$
$\boldsymbol{a}^{(3)} = W^{(3)}\boldsymbol{h}^{(2)}$
$W^{(3)} \in \mathbb{R}^{2\times2}, \boldsymbol{a}^{(3)}, \boldsymbol{h}^{(3)} \in \mathbb{R}^2$

- The number of layers
- The numbers of dimensions of hidden layers
- An activation function for each layer
- A loss function

# Cross entropy loss

- For binary classification

$$l(a, y) = -y \log \sigma(a) - (1 - y) \log\big(1 - \sigma(a)\big)$$

- For multi-class classification

$$l(\boldsymbol{a}, y) = -a_y + \log \sum_k \exp(a_k)$$

- Cross entropy

True probability distribution
(1 for true category; 0 otherwise)

$$H(p, q) = -\sum_k p(k) \log q(k)$$

Predicted probability distribution

# Mean Squared Error (MSE) loss

- Used for regression

$$l(\boldsymbol{a}, \boldsymbol{y}) = \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{a}\|_2^2$$

# Training multi-layer neural networks and back propagation

# Generic notation for multi-layer NNs



First layer: $\mathbb{R}^2 \to \mathbb{R}^3$
$$\boldsymbol{h}^{(1)} = g^{(1)}(\boldsymbol{a}^{(1)})$$
$$\boldsymbol{a}^{(1)} = W^{(1)}\boldsymbol{h}^{(0)}$$
$$W^{(1)} \in \mathbb{R}^{3\times 2}, \boldsymbol{a}^{(1)}, \boldsymbol{h}^{(1)} \in \mathbb{R}^3$$

Second layer: $\mathbb{R}^3 \to \mathbb{R}^2$
$$\boldsymbol{h}^{(2)} = g^{(2)}(\boldsymbol{a}^{(2)})$$
$$\boldsymbol{a}^{(2)} = W^{(2)}\boldsymbol{h}^{(1)}$$
$$W^{(2)} \in \mathbb{R}^{2\times 3}, \boldsymbol{a}^{(2)}, \boldsymbol{h}^{(2)} \in \mathbb{R}^2$$

Final layer: $\mathbb{R}^2 \to \mathbb{R}$
$$\boldsymbol{h}^{(3)} = g^{(3)}(\boldsymbol{a}^{(3)})$$
$$\boldsymbol{a}^{(3)} = W^{(3)}\boldsymbol{h}^{(2)}$$
$$W^{(3)} \in \mathbb{R}^{1\times 2}, \boldsymbol{a}^{(3)}, \boldsymbol{h}^{(3)} \in \mathbb{R}$$

- The $l$–th layer ($l \in \{1, \dots, L\}$) consists of:
  - Input: $\boldsymbol{h}^{(l-1)} \in \mathbb{R}^{d_{l-1}}$      ($\boldsymbol{h}^{(0)} = \boldsymbol{x}$)
  - Output: $\boldsymbol{h}^{(l)} \in \mathbb{R}^{d_l}$      ($\boldsymbol{h}^{(L)} = \widehat{\boldsymbol{y}}$)
  - Weight: $W^{(l)} \in \mathbb{R}^{d_l \times d_{l-1}}$
  - Activation function: $g^{(l)}$
  - Activation: $\boldsymbol{a}^{(l)} \in \mathbb{R}^{d_l}$

$$\boldsymbol{h}^{(l)} = g^{(l)}(W^{(l)}\boldsymbol{h}^{(l-1)})$$

$$W^{(l)} = \left(w_{ij}^{(l)}\right)$$

$w_{ij}^{(l)}$: weight from the $j$-th neuron to the $i$-th neuron of the $l$-th layer

Please accept the notational conflict between an instance-wise loss $l_n$ and a layer number $l$

# How to train weights in MLPs

- We have no explicit supervision signals for the internal (hidden) inputs/outputs $\boldsymbol{h}^{(2)}, \ldots, \boldsymbol{h}^{(L-1)}$

- Having said that, SGD only needs the value of gradient $\frac{\partial l_n}{\partial w_{ij}^{(l)}}$ for every weight $w_{ij}^{(l)}$ in MLPs

- Can we compute the value of $\frac{\partial l_n}{\partial w_{ij}^{(l)}}$ for every weight $w_{ij}^{(l)}$?
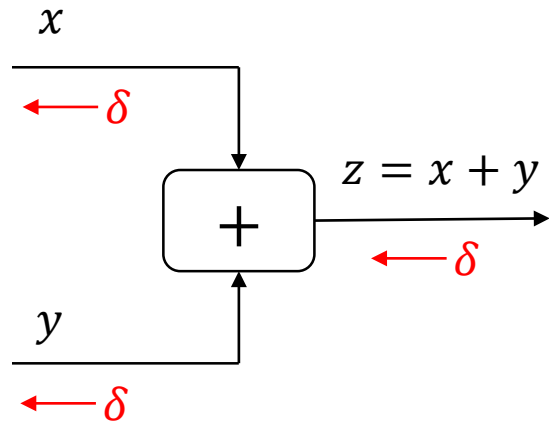
- Yes! *Backpropagation* can do that!!

# Backpropagation

- Commonly used in deep neural networks

- Formulas for backpropagation look complicated

- However:
  - We can understand backpropagation easily if we know the concept of *computation graph*
  - Most deep learning frameworks implement backpropagation by using *automatic differentiation*

- Let's see computation graph and automatic differentiation first

## General Back-Propagation

The back-propagation algorithm is very simple. To compute the gradient of some scalar $z$ with respect to one of its ancestors $\boldsymbol{x}$ in the graph, we begin by observing that the gradient with respect to $z$ is given by $\frac{dz}{dz} = 1$. We can then compute the gradient with respect to each parent of $z$ in the graph by multiplying the current gradient by the Jacobian of the operation that produced $z$. We continue multiplying by Jacobians traveling backwards through the graph in this way until we reach $\boldsymbol{x}$. For any node that may be reached by going backwards from $z$ through two or more paths, we simply sum the gradients arriving from different paths at that node.
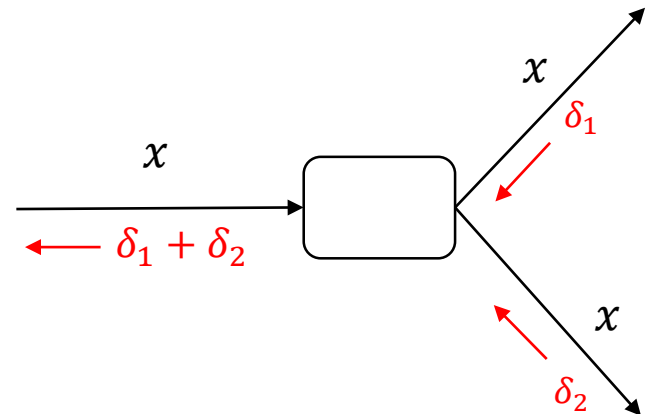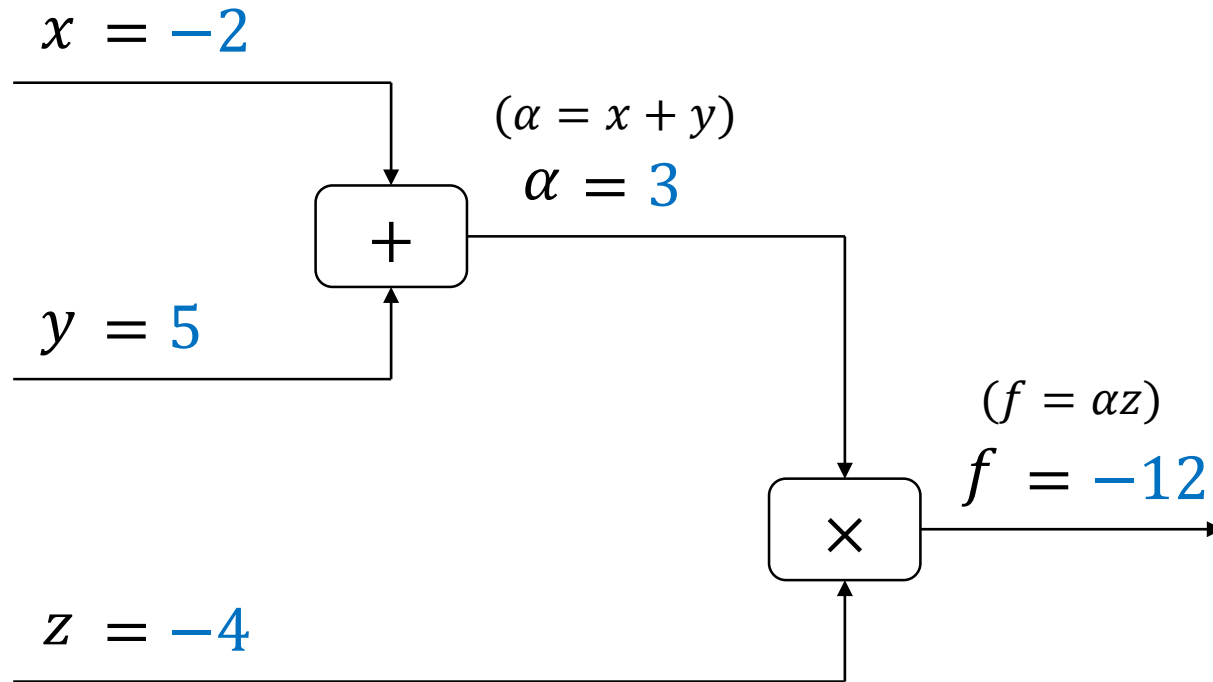
# Rules for reverse-mode AD

$x$

$\longleftarrow \delta$

$$z = x + y$$

$+$

$\longleftarrow \delta$

$y$

$\longleftarrow \delta$

Add

$x$ $\qquad z = f(x)$

$\dfrac{\partial f(x)}{\partial x} \cdot \delta$ $\qquad f(x)$

$\longleftarrow$ $\qquad\qquad \longleftarrow \delta$

Function application

$x$

$\longleftarrow y \cdot \delta$

$$z = xy$$

$\times$

$\longleftarrow \delta$

$y$

$\longleftarrow x \cdot \delta$

Multiply

$x$

$\delta_1$

$x$

$\longleftarrow \delta_1 + \delta_2$

$x$

$\delta_2$

Branch

# Computation graph: $f(x, y, z) = (x + y)z$

$x = -2$

$(\alpha = x + y)$

$\alpha = 3$

$y = 5$

$+$

$(f = \alpha z)$

$f = -12$

$\times$

$z = -4$

*Forward pass*

The value of a variable (above an arrow)

# Automatic Differentiation (AD): $f(x, y, z) = (x + y)z$

http://cs231n.github.io/optimization-2/

$x = -2$

$\longleftarrow -4$

$\left( \dfrac{\partial \alpha}{\partial x} \times (-4) \right)$

$(\alpha = x + y)$

$\alpha = 3$

$+$

$\longleftarrow -4$

$\left( \dfrac{\partial f}{\partial \alpha} \times 1 = z = -4 \right)$

$y = 5$

$\longleftarrow -4$

$\left( \dfrac{\partial \alpha}{\partial y} \times (-4) \right)$

$(f = \alpha z)$

$f = -12$

$\times$

$\longleftarrow 1$

$z = -4$

$\longleftarrow 3$

$\left( \dfrac{\partial f}{\partial z} \times 1 = \alpha = 3 \right)$

*Backward pass*
*(Reverse mode AD)*

Compare with:

$\dfrac{\partial f}{\partial x} = z = -4$

$\dfrac{\partial f}{\partial y} = z = -4$

$\dfrac{\partial f}{\partial z} = (x + y) = 3$

The value of a variable (above an arrow)
The gradient of the output $f$ with respect to the variable (below an arrow)

# Automatic differentiation (Baydin+ 2018)

- AD computes derivations by using the *chain rule*
  - Function values computed in the forward pass
  - Derivations computed with respect to:
    - Every variable (in reverse-mode accumulation)
    - A specific variable (in forward-mode accumulation)

- Do not confuse with these:
  - Numerical differentiation: for example, $\frac{\partial f(x)}{\partial x} = \frac{f(x+\delta) - f(x)}{\delta}$
  - Symbolic differentiation: e.g., Mathematica, sympy

# Exercise: AD on computation graph

- Write a computation graph for $l_x(w)$,

$$l_x(w) = -\log \sigma(\boldsymbol{w} \cdot \boldsymbol{x}) = -\log \frac{1}{1 + e^{-\boldsymbol{w} \cdot \boldsymbol{x}}}$$

- Consider $\boldsymbol{x} = (1,1,1)^\top$ and $\boldsymbol{w} = (1,1,-1.5)^\top$
  - Compute the value of $l_x(\boldsymbol{w})$
  - Compute gradients $\frac{\partial l_x(\boldsymbol{w})}{\partial \boldsymbol{w}}$

# Computing $\frac{\partial l_x(\boldsymbol{w})}{\partial \boldsymbol{w}}$ using AD



$$\gamma = \alpha\beta \qquad \partial\gamma/\partial\alpha = \beta \qquad \partial\gamma/\partial\beta = \alpha$$

$$\zeta = \delta\varepsilon \qquad \partial\zeta/\partial\delta = \varepsilon \qquad \partial\zeta/\partial\varepsilon = \delta$$

$$\kappa = \theta\vartheta \qquad \partial\kappa/\partial\theta = \vartheta \qquad \partial\kappa/\partial\vartheta = \theta$$

$$\lambda = \gamma + \zeta \qquad \partial\lambda/\partial\gamma = 1 \qquad \partial\lambda/\partial\zeta = 1$$

$$\mu = \lambda + \kappa \qquad \partial\mu/\partial\lambda = 1 \qquad \partial\mu/\partial\kappa = 1$$

$$\nu = -\mu \qquad \partial\nu/\partial\mu = -1$$

$$\xi = e^\nu \qquad \partial\xi/\partial\nu = e^\nu$$

$$\pi = \xi + 1 \qquad \partial\pi/\partial\xi = 1$$

$$\varpi = 1/\pi \qquad \partial\varpi/\partial\pi = -(1/\pi)^2$$

$$\rho = \log\varpi \qquad \partial\rho/\partial\varpi = 1/\varpi$$

$$l = -\rho \qquad \partial l/\partial\rho = -1$$

$$w_1 \leftarrow w_1 + \eta\frac{\partial l_x(\boldsymbol{w})}{\partial w_1} = w_1 + 0.3775\eta$$

$$\frac{\partial\rho}{\partial\varpi} \times (-1)$$
$$= -\frac{1}{\varpi} = -\frac{1}{0.6225}$$
$$= -1.6065$$

$$\frac{\partial\varpi}{\partial\pi} \times (-1.6065)$$
$$= -\left(\frac{1}{1.6065}\right)^2 \times (-1.6065)$$
$$= 0.6224$$

# No need to derive backpropagation

- Manual derivation of gradients is tedious and error-prone
  - Debugging a mistake in gradients is extremely difficult

- AD is employed in most deep learning frameworks
  - We only need implement an algorithm for a forward pass, i.e., how to compute an output from an input
  - We can concentrate on designing a structure of neural network
  - This boosted the speed of research and development
  - The idea of AD is not new (since 1959)

- Deriving a formula for backpropagation is legacy

# Summary and notes

- We design:
  - A neural network model $f(\boldsymbol{x}; \theta)$ (with parameters $\theta$)
  - A loss function: $E_D(\theta) = \sum_{n=1}^{N} \mathcal{L}(f(\boldsymbol{x}_n; \theta), y_n)$
    - $\mathcal{L}$ is an instance-wise loss function
    - $D$ presents a set of training data $D = \left((x_1, y_1), \ldots, (x_N, y_N)\right)$

- We find a minimizer $\theta^*$ for $E_D(\theta)$ by using SGD
  - An update formula for every parameter $w \in \theta$ is derived in a generic manner based on automatic differentiation

- Step function is inappropriate for backpropagation
  - Gradients will not flow because $g'(a) = 0$ at $a \neq 0$

# Activation functions

# Step

Step function: $\mathbb{R} \to \{0,1\}$

$$g(x) = \begin{cases} 1 & (\text{if } x > 0) \\ 0 & (\text{otherwise}) \end{cases}$$



- Pros
  - Yields a binary output
- Cons (never use this)
  - Zero gradients
    - SGD cannot update parameters because $\frac{\partial l}{\partial w} = 0$

# Sigmoid

Sigmoid: $\mathbb{R} \to (0,1)$

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



- Pros
  - Yields an output within (0,1)
- Cons
  - Not zero-centered
  - Zero (vanishing) gradients when $|x|$ is large

# Hyperbolic tangent (tanh)

$$\tanh: \mathbb{R} \to (-1,1)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\sigma(2x) - 1$$



- Pros
  - Yields an output within $(-1,1)$
  - Zero-centered
- Cons
  - Zero (vanishing) gradients when $|x|$ is large

# Rectified Linear Unit (ReLU)

$$\text{ReLU}: \mathbb{R} \to \mathbb{R}_{\geq 0}$$
$$\text{ReLU}(x) = \max(0, x)$$



- Pros
  - Gradients do not vanish when $x > 0$
  - Light-weight (no $e^x$) computation
  - Faster convergence (e.g., 6x faster on CIFAR-10)
- Cons
  - Not zero centered
  - Dead neurons when $x \leq 0$

# Leaky ReLU

Leaky ReLU: $\mathbb{R} \to \mathbb{R}$
$$\text{LeakyReLU}_\alpha(x) = \max(\alpha x, x)$$



- Pros
  - Gradients do not vanish
  - Light-weight (no $e^x$) computation
- Cons
  - Not zero centered
  - Not so much improvement over ReLU in practice

# Typical definition of a DNN

Typical example of a (pseudo) definition for a Deep Neural Network of depth L (L-1 hidden layers and 1 output layer) for (almost) all modern DL frameworks and libraries:

```
initialize dnnmodel
```

dataset consists of D inputs x of dimensions d each

set hidden Layer 1 as a fully connected (fc) layer to inputs x containing n1 neurons (n1 x d connections):

```
dnnmodel.fc1(n1)
fc1.activation-function=relu
```

set hidden Layer 2 as a fully connected layer to Layer 1 containing n2 neurons (n2 x n1 connections):

```
dnnmodel.fc2(n2)
fc2.activation-function=relu
```

… more hidden layers ...

set hidden Layer L-1 as a fully connected layer to Layer L-2 containing nL neurons (nL-1 x nL-2 connections):

```
dnnmodel.fcL-1(nL-1)
fcL-1.activation-function=relu
```

set fcL (output layer) in a binary classification problem (1 x nL-1 connections)

```
dnnmodel.fcL(1)
fcL.activation-function=sigmoid
```

OR set fcL (output layer) in a K-class multiclass classification problem with one-hot encoding (K x nL-1 connections)

```
dnnmodel.fcL(K)
fcL.activation-function=softmax
```

Set Loss function to Cross Entropy (based on Maximum Likelihood Estimation). Use C for L2 weight regularization.

```
dnnmodel.loss=crossentropy(C)
```

Set the solver to mini batch SGD (or variations) with learning rate lr that by default will use automatic differentiation for the backprop gradient computation. Also set maximum number of epochs (e) and batch size (m) and use of dropout for further regularization. (one epoch equals D examples used for training)

```
dnnmodel.solver=minibatchSGD(lr,m,e,dropout=yes)
```

# What made Deep Neural Networks possible and efficient?

Factors from the natural evolution of computation
- Better computers and software allow bigger networks with higher capacity to solve more difficult problems
- With bigger datasets available we must use stochastic methods like SGD

Algorithmic factors
- Cross-entropy is a better loss function than MSE for sigmoid, softmax
- ReLU in hidden layers is a better activation function than sigmoid and tanh for deeper networks
- Automatic differentiation is now a feature of all DL frameworks

# What is the difference between a neural network and a deep neural network, and why do the deep ones work better?

Short answer: DNN simply seem to perform better! Read the first answer in the following link:

https://stats.stackexchange.com/questions/182734/what-is-the-difference-between-a-neural-network-and-a-deep-neural-network-and-w