# CS 6501: Deep Learning for Computer Graphics

# **Basics of Neural Networks**

**Connelly Barnes** 

#### Overview

- Simple neural networks
  - Perceptron
  - Feedforward neural networks
  - Multilayer perceptron and properties
  - Autoencoders
- How to train neural networks
  - Gradient descent
  - Stochastic gradient descent
  - Automatic differentiation
  - Backpropagation

### Perceptron (1957, Cornell)



# Perceptron (1957, Cornell)

• Binary classifier, can learn linearly separable patterns.



Diagram from Wikipedia

#### Feedforward neural networks

• We could connect units (neurons) in any arbitrary graph



## Feedforward neural networks

- We could connect units (neurons) in any arbitrary graph
- If no cycles in the graph we call it a **feedforward neural network**.



### Recurrent neural networks (later)

• If cycles in the graph we call it a **recurrent neural network**.



#### Overview

- Simple neural networks
  - Perceptron
  - Feedforward neural networks
  - Multilayer perceptron and properties
  - Autoencoders
- How to train neural networks
  - Gradient descent
  - Stochastic gradient descent
  - Automatic differentiation
  - Backpropagation

# Multilayer Perceptron (1960s)



In matrix notation:  $L_{i} = f_{i}(W_{i}L_{i-1} + b_{i}), i \ge 1$ L<sub>1</sub>: Input layer (inputs) vector L<sub>2</sub>: Hidden layer vector L<sub>3</sub>: Output layer vector  $W_i$ : Weight matrix for connections from layer *i*-1 to layer *i* 

- $\mathbf{b}_i$ : Biases for neurons in layer *i* 
  - $f_i$ : Activation function for layer *i*

#### Activation Functions: Sigmoid / Logistic



#### **Activation Functions: Tanh**



# Activation Functions: ReLU (Rectified Linear Unit)

$$f(x) = \max(x, 0)$$

 $\frac{df}{dx} = \begin{cases} 1, & \text{if } x > 0\\ 0, & \text{if } x < 0 \end{cases}$ 

Pros:

- Accelerates training stage by 6x over sigmoid/tanh [1]
- Simple to compute
- Sparser activation patterns

#### Cons:

• Neurons can "die" by getting stuck in zero gradient region Summary:

2

• Currently preferred kind of neuron

## **Universal Approximation Theorem**

- Multilayer perceptron with a single hidden layer and linear output layer can approximate any continuous function on a compact subset of R<sup>n</sup> to within any desired degree of accuracy.
- Assumes activation function is bounded, non-constant, monotonically increasing.
- Also applies for <u>ReLU activation function</u>.

# **Universal Approximation Theorem**

- In the worst case, exponential number of hidden units may be required.
- Can informally show this for binary case:
  - If we have *n* bits input to a binary function, how many possible inputs are there?
  - How many possible binary functions are there?
  - So how many weights do we need to represent a given binary function?

# Why Use Deep Networks?

- Functions representable with a deep rectifier network can require an exponential number of hidden units with a shallow (one hidden layer) network (Goodfellow 6.4)
- Piecewise linear networks (e.g. using ReLU) can represent functions that have a number of regions exponential in depth of network.
  - Can capture repeating / mirroring / symmetric patterns in data.
  - Empirically, greater depth often results in better generalization.

#### Neural Network Architecture

- Architecture: refers to which parameters (e.g. weights) are used in the network and their topological connectivity.
- **Fully connected:** A common connectivity pattern for multilayer perceptrons. All possible connections made between layers *i*-1 and *i*.



Is this network fully connected?

## Neural Network Architecture

- Architecture: refers to which parameters (e.g. weights) are used in the network and their topological connectivity.
- **Fully connected:** A common connectivity pattern for multilayer perceptrons. All possible connections made between layers *i*-1 and *i*.



Is this network fully connected?

# How to Choose Network Architecture?

- Long discussion
- Summary:
  - Rules of thumb do not work.
    - "Need 10x [or 30x] more training data than weights."
    - Not true if very low noise
    - Might need even more training data if high noise
  - Try many networks with different numbers of units and layers
  - Check generalization using validation dataset or cross-validation.

#### Overview

- Simple neural networks
  - Perceptron
  - Feedforward neural networks
  - Multilayer perceptron and properties
  - Autoencoders
- How to train neural networks
  - Gradient descent
  - Stochastic gradient descent
  - Automatic differentiation
  - Backpropagation

### Autoencoders

- Learn the identity function  $h_w(\mathbf{x}) = \mathbf{x}$
- Is this supervised or unsupervised learning?



# Autoencoders

- Applications:
  - Dimensionality reduction
  - Learning manifolds
  - Hashing for search problems







#### Overview

- Simple neural networks
  - Perceptron
  - Feedforward neural networks
  - Multilayer perceptron and properties
  - Autoencoders
- How to train neural networks
  - Gradient descent
  - Stochastic gradient descent
  - Automatic differentiation
  - Backpropagation

#### **Gradient Descent**

Discuss amongst students near you:

- What are some problems that could be easily optimized with gradient descent?
- Problems where this is difficult?
- Should the learning rate be constant or change?



# Gradient Descent with Energy Functions that have Narrow Valleys



Source: "Banana-SteepDesc" by P.A. Simionescu – Wikipedia English

#### **Gradient Descent with Momentum**

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \Delta_n$$
$$\Delta_n = \gamma_n \nabla F(\mathbf{x}_n) + m \Delta_{n-1}$$
Momentum

Could use small value e.g. m=0.5 at first Could use larger value e.g. m=0.9 near end of training when there are more oscillations.

#### **Gradient Descent with Momentum**



#### Without Momentum

With Momentum

Figure from Genevieve B. Orr, Willamette.edu

# Stochastic gradient descent

- <u>Stochastic gradient descent</u> (Wikipedia)
- Gradient of sum of n terms where n is large
- Sample rather than computing the full sum
  - Sample size s is "mini-batch size"
  - Could be 1 (very noisy gradient estimate)
  - Could be 100 (collect photos 100 at a time to find each noisy "next" estimate for the gradient)
- Use same step as in gradient descent to the estimated gradient

#### Stochastic gradient descent

• Pseudocode:

- Repeat until an approximate minimum is obtained:
  - Randomly shuffle examples in the training set.
  - For  $i=1,2,\ldots,n$  , do:

 $ullet w := w - \eta 
abla Q_i(w).$ 



#### **Problem Statement**

- Take the gradient of an arbitrary program or model (e.g. a neural network) with respect to the parameters in the model (e.g. weights).
- If we can do this, we can use gradient descent!

#### **Review: Chain Rule in One Dimension**

- Suppose  $f: \mathbb{R} \to \mathbb{R}$  and  $g: \mathbb{R} \to \mathbb{R}$
- Define

$$h(x) = f(g(x))$$

• Then what is h'(x) = dh/dx?

$$h'(x) = f'(g(x))g'(x)$$

# Chain Rule in Multiple Dimensions

- Suppose  $f: \mathbb{R}^m \to \mathbb{R}$  and  $g: \mathbb{R}^n \to \mathbb{R}^m$ , and  $\mathbf{x} \in \mathbb{R}^n$
- Define

$$h(\mathbf{x}) = f(g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))$$

 Then we can define partial derivatives using the <u>multidimensional</u> <u>chain rule</u>:

$$\frac{\partial f}{\partial x_i} = \sum_{l=1}^m \frac{\partial f}{\partial g_j} \frac{\partial g_j}{\partial x_i}$$

# The Problem with Symbolic Derivatives

• What if our program takes 5 exponential operations:

$$y = \exp\left(\exp\left(\exp\left(\exp\left(\exp(\exp(x)\right)\right)\right)\right)$$

- What is *dy/dx*? (blackboard)
- How many exponential operations in the resulting expression?
- What if the program contained *n* exponential operations?

# Solution: Automatic Differentiation (1960s, 1970s)

- Write an arbitrary program as consisting of basic operations f<sub>1</sub>, ..., f<sub>n</sub>
   (e.g. +, -, \*, cos, sin, ...) that we know how to differentiate.
- Label the inputs of the program as  $x_1, \ldots, x_n$ , and the output  $x_N$ .
- The computation:

For i = n + 1, ..., N  $x_i = f_i(x_{\pi(i)})$   $\pi(i)$ : Sequence of "parent" values (e.g. if  $\pi(3) = (1,2)$ , and  $f_3 = +$ , then  $x_3 = x_1 + x_2$ )

• Reverse mode automatic differentiation: apply the chain rule from the end of the program  $x_N$  back towards the beginning.

# Solution for Simplified Chain of Dependencies

- Suppose  $\pi(i) = i 1$
- The computation: For i = n + 1, ..., N $x_i = f_i(x_{i-1})$



For example:

• What is 
$$\frac{dx_N}{dx_N}$$
 ?

# Solution for Simplified Chain of Dependencies

- Suppose  $\pi(i) = i 1$
- The computation: For i = n + 1, ..., N

 $x_i = f_i(x_{i-1})$ 



For example:

$$\frac{dx_N}{dx_N} = 1$$

#### Solution for Simplified Chain of Dependencies For example:

- St
- Th

• Suppose 
$$\pi(i) = i - 1$$
  
• The computation:  
For  $i = n + 1, ..., N$   
 $x_i = f_i(x_{i-1})$ 
• What is  $\frac{dx_N}{dx_i}$  in terms of  $\frac{dx_N}{dx_{i+1}}$ ?  
 $dx_{i+1} = dx_i \left(\frac{\partial x_{i+1}}{\partial x_i}\right) \quad \Box \gg \frac{dx_{i+1}}{dx_N} = \frac{dx_i}{dx_N} \left(\frac{\partial x_{i+1}}{\partial x_i}\right)$ 

# Solution for Simplified Chain of Dependencies

- Suppose  $\pi(i) = i 1$
- The computation: For i = n + 1, ..., N $x_i = f_i(x_{i-1})$



- What is  $\frac{dx_N}{dx_i}$  in terms of  $\frac{dx_N}{dx_{i+1}}$ ?  $\frac{dx_N}{dx_i} = \frac{dx_N}{dx_{i+1}} \left(\frac{\partial x_{i+1}}{\partial x_i}\right)$
- Conclusion: run the computation forwards. Then initialize  $\frac{dx_N}{dx_N} = 1$ and work **backwards** through the computation to find  $\frac{dx_N}{dx_i}$  for each *i* from  $\frac{dx_N}{dx_{i+1}}$ . This gives us the gradient of the output  $(x_N)$  with respect to every expression in our compute graph!

# What if the Dependency Graph is More Complex?



• The computation:

For i = n + 1, ..., N $x_i = f_i(\mathbf{x}_{\pi(i)})$ 

 $\pi(i)$ : Sequence of "parent" values (e.g. if  $\pi(3) = (1,2)$ , and  $f_3 = +$ , then  $x_3 = x_1 + x_2$ )

Solution: apply multi-dimensional chain rule.

## Solution: Automatic Differentiation (1960s, 1970s)

• Computation:

$$x_i = f_i(\mathbf{x}_{\pi(i)})$$

• Multidimensional chain rule:

$$\frac{\partial f}{\partial x_i} = \sum_{l=1}^m \frac{\partial f}{\partial g_j} \frac{\partial g_j}{\partial x_i}$$

• Result:

$$\frac{dx_N}{dx_i} = \sum_{j:i\in\pi(j)} \frac{dx_N}{dx_j} \frac{\partial x_j}{\partial x_i}$$

 $f(g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))$ 

#### Explanation from Justin Domke

# Solution: Automatic Differentiation (1960s, 1970s)

• Back-propagation algorithm: initialize:

$$\frac{dx_N}{dx_N} = 1$$

Example on blackboard for a program with one addition.

$$\frac{dx_N}{dx_i} = \sum_{j:i\in\pi(j)} \frac{dx_N}{dx_j} \frac{\partial x_j}{\partial x_i}$$

• Now we have differentiated the output of the program  $x_N$  with respect to the inputs  $x_1, \dots, x_n$ , as well as every computed expression  $x_i$ .

Explanation from Justin Domke

- Apply reverse mode automatic differentiation to a neural network's loss function.
- A special case of what we just derived.
- If we have one output neuron, squared error is:

$$E=rac{1}{2}(t-y)^2$$

t is the target output for a training sample, and

y is the actual output of the output neuron.



For each neuron j, its output  $o_j$  is defined as

$$o_j = arphi(\mathrm{net}_j) = arphi\left(\sum_{k=1}^n w_{kj}o_k
ight).$$

The input  $net_j$  to a neuron is the weighted sum of outputs  $o_k$  of previous neurons. If The variable  $w_{ij}$  denotes the weight between neurons i and j.



From Wikipedia

• Apply the chain rule twice:

$$rac{\partial E}{\partial w_{ij}} = rac{\partial E}{\partial o_j} rac{\partial o_j}{\partial \mathrm{net_j}} rac{\partial \mathrm{net_j}}{\partial w_{ij}}$$

• Last term is easy:

$$egin{aligned} & o_j = arphi( ext{net}_j) = arphi\left(\sum_{k=1}^n w_{kj}o_k
ight). \ & rac{\partial ext{net}_j}{\partial w_{ij}} = rac{\partial}{\partial w_{ij}}\left(\sum_{k=1}^n w_{kj}o_k
ight) = o_i \end{aligned}$$





• Apply the chain rule twice:

$$rac{\partial E}{\partial w_{ij}} = rac{\partial E}{\partial o_j} rac{\partial o_j}{\partial \mathrm{net_j}} rac{\partial \mathrm{net_j}}{\partial w_{ij}}$$

• Second term is easy:

$$o_j = \varphi(\operatorname{net}_j)$$
  
 $rac{\partial o_j}{\partial net_j} = \varphi'(\operatorname{net}_j)$ 





• Apply the chain rule twice:

$$rac{\partial E}{\partial w_{ij}} = rac{\partial E}{\partial o_j} rac{\partial o_j}{\partial \mathrm{net_j}} rac{\partial \mathrm{net_j}}{\partial w_{ij}}$$

$$egin{aligned} & o_j = y \ & rac{\partial E}{\partial o_j} = \end{aligned}$$
 (Derivation on board)

$$E=rac{1}{2}(t-y)^2$$

.



From Wikipedia

• Apply the chain rule twice:

$$rac{\partial E}{\partial w_{ij}} = rac{\partial E}{\partial o_j} rac{\partial o_j}{\partial \mathrm{net_j}} rac{\partial \mathrm{net_j}}{\partial w_{ij}}$$

- If the neuron is interior neuron, we use the chain rule from automatic differentiation.
- To do this, we need to know what expressions depend on the current neuron's output o<sub>j</sub>?
- Answer: other neurons input sums, i.e. net<sub>l</sub> for all neurons *l* receiving inputs from the current neuron.





• Apply the chain rule twice:

$$rac{\partial E}{\partial w_{ij}} = rac{\partial E}{\partial o_j} rac{\partial o_j}{\partial \mathrm{net_j}} rac{\partial \mathrm{net_j}}{\partial w_{ij}}$$



• If the neuron is an interior neuron, chain rule:

$$rac{\partial E}{\partial o_j} = \sum_{l \in L} \left( rac{\partial E}{\partial \mathrm{net}_l} rac{\partial \mathrm{net}_l}{\partial o_j} 
ight) \ = \sum_{l \in L} \left( rac{\partial E}{\partial o_l} rac{\partial o_l}{\partial \mathrm{net}_l} w_{jl} 
ight)$$

All neurons receiving input from the current neuron.

From Wikipedia

• Partial derivative of error *E* with respect to weight  $w_{ij}$ :

$$\frac{\partial E}{\partial w_{ij}} = \delta_j o_i$$
  

$$\delta_j = \frac{\partial E}{\partial o_j} \frac{\partial o_j}{\partial \operatorname{net}_j} = \varphi'(o_j) \begin{cases} (o_j - t_j) & \text{if } j \text{ is an output neuron} \\ \sum_{l \in L} \delta_l w_{jl} & \text{if } j \text{ is an interior neuron} \end{cases}$$

All neurons receiving input from the current neuron.



Forward direction



• Calculate network and error.



• Backpropagate: from output to input, recursively compute  $\frac{\partial W_{ii}}{\partial W_{ii}} = \nabla_W E$ 

# Gradient Descent with Backpropagation

- Initialize weights at good starting point w<sub>0</sub>
- Repeatedly apply gradient descent step (1)
- Continue training until validation error hits a minimum.



# Stochastic Gradient Descent with Backpropagation

X<sub>2</sub>

- Initialize weights at good starting point w<sub>0</sub>
- Repeat until validation error hits a minimum:
  - Randomly shuffle dataset
  - Loop through mini-batches of data, batch index is *i*
  - Calculate stochastic gradient using backpropagation for each, and apply update rule (1)



 $\mathbf{w}_{n+1} = \mathbf{w}_n - \gamma_n \nabla_w E_i (\mathbf{w}_n)$