### ECE 6775 High-Level Digital Design Automation Fall 2024

# **Introduction to Neural Networks**

Jordan Dotzel, Ritchie Zhao, Zhiru Zhang School of Electrical and Computer Engineering



**Cornell University** 



### Announcements

- Lecture tomorrow 4:30pm at Rhodes 310
- No Friday TA OH this week
- TA OH today 5pm ~ 6pm
  - In-person at Rhodes 312
- No instructor OH today
- Lab 3 due next Monday

### **Rise of Deep Neural Networks (DNNs)**

- DNNs have revolutionized information technology
  - computer vision, e-commerce, finance, game AI, healthcare, machine transcription & translation, robots, web search, and many more (to come)



### **A Brief History**

- 1940's 1950's: First artificial neural networks proposed based on biological structures in the human visual cortex
- 1980's 2000's: Neural networks (NNs) considered inferior to other simpler algorithms (e.g., SVM)
- Mid 2000's: NN research considered "dead", machine learning conferences outright reject most NN papers
- 2010 2012: DNNs begin winning large-scale image and document classification contests
- 2012 Now: DNNs prove themselves in many industrial applications (web search, translation, image analysis)

### Neural Network (NN) in a Nutshell

- Depends on a large volume of data & mostly supervised
- Learns a function that encodes a hierarchy of abstract features
  - Consists of a stack of connected *layers*, such as convolutional, pooling, full connected, attention



### **NN Training and Inference**

- Training refers to the process of building an NN model to accomplish a specific AI task by "learning" from a predetermined dataset
- Inference refers to the use of a trained NN model to make a prediction (or decision) on unseen data



Part 1

# CLASSIFICATION WITH THE PERCEPTRON

### **Classification Problems**

We'll discuss neural networks for solving supervised classification problems



- Given a training set consisting of labeled inputs
- Learn a function which maps inputs to labels
- This function is used to predict the labels of new inputs



**Predict New Data** 

### **Artificial Neuron**

The simplest possible neural network contains only one "neuron", which is described by the following equation:

$$y = \sigma(\sum_{i=1}^{n} w_i x_i + b)$$
  
$$w_i = \text{weights}$$
  
$$b = \text{bias}$$
  
$$\sigma = \text{activation function}$$

- When σ is the unit step (or sign) function, we have a perceptron\*
  - Invented in 1957 by Frank Rosenblatt
  - \* Perceptron is often used as a synonym for artificial neuron, where  $\sigma$  could be any activation function

### **Activation Function**

#### • The activation function $\sigma$ is non-linear



Used in the initial perceptrons



neural nets

Makes deep networks easier to train

0.5

Used in modern deep nets

#### **Breaking Down the "Neuron"**



#### **Breaking Down the "Neuron"**

A 2-input, 1-output neuron



A real-life analogue

#### **A Real-life Classification Problem**

- Inputs: Pairs of numbers  $(x_1, x_2)$
- Labels: 0 or 1 (binary decision problem)
- Real-life analogue:
  - Label = Raining or Not Raining
  - $x_1$  = Relative humidity
  - $x_2$  = Cloud coverage

<b>x</b> <sub>1</sub>	<b>X</b> 2	Label
-0.7	-0.6	0
-0.6	0.4	0
-0.5	0.7	1
-0.5	-0.2	0
-0.1	0.7	1
0.1	-0.2	0
0.3	0.5	1
0.4	0.1	1
0.4	-0.7	0
0.7	0.2	1

**Training set** 

#### **Visualizing the Data**



<b>X</b> 1	<i>X</i> <sub>2</sub>	Label	
-0.7	-0.6	0	
-0.6	0.4	0	
-0.5	0.7	1	
-0.5	-0.2	0	
-0.1	0.7	1	
0.1	-0.2	0	
0.3	0.5	1	
0.4	0.1	1	
0.4	-0.7	0	
0.7	0.2	1	
Training set			

13

#### **Decision Boundary**



In this case, the data points can be classified with a **linear decision boundary** produced by a perceptron

<i>x</i> <sub>1</sub>	<b>X</b> 2	Label
-0.7	-0.6	0
-0.6	0.4	0
-0.5	0.7	1
-0.5	-0.2	0
-0.1	0.7	1
0.1	-0.2	0
0.3	0.5	1
0.4	0.1	1
0.4	-0.7	0
0.7	0.2	1

#### **Training set**

### **Finding the Parameters**

- The right parameters (weights and bias) will create any linear decision boundary we want
- Training = process of finding the parameters to solve our classification problem
  - Basic idea: iteratively modify the parameters to reduce the training loss
  - Training loss: measure of difference between predictions and labels on the training set

#### **Gradient Descent**

#### Loss function

Measure of difference between predictions and true labels

$$L = \sum_{\substack{i=0 \\ \uparrow}}^{N} (y^{(i)} - t^{(i)})^2$$
  
Sum over training samples

$$y^{(i)}$$
 = Prediction  
 $t^{(i)}$  = True label

Over dia set diversions

Gradient Descent:

$$w_{k+1} = w_k - \eta \frac{\partial L}{\partial w_k} \leftarrow \begin{array}{c} \text{Gradient = direction} \\ \text{of steepest descent} \\ \text{in } L \end{array}$$

k = training step  $\eta$  = learning rate or step size

### **Training a Neural Network**

- At each step k:
  - 1. Classify each sample to get each  $y^{(i)}$
  - 2. Compute the **loss** *L*
  - 3. Compute the gradient  $\frac{\partial L}{\partial w_k}$
  - 4. Update the parameters using **gradient descent**  $w_{k+1} = w_k - \eta \frac{\partial L}{\partial w_k}$

#### Demo

### Perceptron training demo

- No bias (bias = 0)
- No test set (training samples only)

#### **Another Example**

<i>x</i> <sub>2</sub>	<i>x</i> <sub>1</sub>	OR		
0	0	0	$w_1 x_1 + w_2 x_2 + b < 0$	$\Rightarrow b < 0$
0	1	1	$w_1 x_1 + w_2 x_2 + b \ge 0$	$\implies w_1 \ge -b$
1	0	1	$w_1 x_1 + w_2 x_2 + b \ge 0$	$\implies w_2 \ge -b$
1	1	1	$w_1 x_1 + w_2 x_2 + b < 0$	$\implies w_1 + w_2 \ge -b$



One solution:

$$w_1 = 1; w_2 = 1; b = -1$$

19

#### Yet Another Example (The XOR Problem)

<i>x</i> <sub>2</sub>	$x_1$	XOR		
0	0	<mark>0</mark> W	$_1x_1 + w_2x_2 + b < 0$	$\Rightarrow b < 0$
0	1	1 W	$_{1}x_{1} + w_{2}x_{2} + b \ge 0$	$\implies w_1 \ge -b$
1	0	1 W	$_1x_1 + w_2x_2 + b \ge 0$	$\implies w_2 \ge -b$
1	1	<mark>0</mark> W	$_{1}x_{1} + w_{2}x_{2} + b < 0$	$\implies w_1 + w_2 < -b$



## Part 2 DEEP NEURAL NETWORKS

#### **Combining Neurons**

- A single neuron can only make a simple decision
- Feeding neurons into each other allows a neural network to learn complex decision boundaries



Source: http://www.opennn.net/

#### **Complex Decision Boundaries**



Source: https://www.carl-olsson.com/fall-semester-2013/

#### **Multi-Layer Perceptron (MLP)**

- MLP is a fully connected class of feedforward artificial neural network (ANN)
  - An MLP model consists of multiple layers of neurons
  - Often referred to as "vanilla" ANNs, especially with a single hidden layer



A vectorized (tensorized) view of MLP



#### **Deep Neural Network (DNN)**

MLPs are neural networks with at least three layers, while DNNs typically have even more (hidden) layers



Image credit: http://www.opennn.net/

#### Learning a Deep Neural Network

Gradient Descent:

$$w_{k+1} = w_k - \eta \frac{\partial L}{\partial w_k}$$



#### **Backpropagation**

Backpropagation: use the chain rule from calculus to propagate the gradients backwards through the network



#### **Stochastic Gradient Descent**

#### Remember Gradient Descent? $\partial L$ $w_{k+1} = w_k - \eta \overline{\partial w_k}$

ich can be millions of samples!

#### Stochastic Gradient Descent:

- At each set, only compute *L* for a **minibatch** (a few samples randomly taken from the training set)
- SGD is faster and more accurate than GD for DNNs!

Part 3 CONVOLUTIONAL NEURAL NETWORKS

### **Neural Networks for Images**

- So far, we've seen neural networks built from fully-connected layers
- Do we really need all the edges for learning an image?
  - Important patterns are typically much smaller than the whole image
  - Images are also shift-invariant (e.g., a bird is a bird even when shifted)

input layer a hidden layer 1 hidden layer 2 hidden layer 3





### **Convolutional Neural Network (CNN)**



- Front: convolutional layers learn visual features
- Feature maps get downsampled through the network
- Back: fully-connected layers perform classification using the visual features

### **A Convolutional (conv) Layer**

- A CNN stacks multiple conv layers (and some other layers)
- A conv layer has a set of learnable filters that perform convolution operation



### **The Convolutional Filter**



- Each neuron learns a weight filter and convolves the filter over the image
- Each neuron outputs a 2D feature map (basically an image of features)

### **The Convolutional Filter**



- Each point in the feature map encodes both a decision and its spatial location
- Detects the pattern anywhere in the image!

#### **The Convolutional Layer: A Closer Look**



- M input and N output feature maps
- Each output map uses *M* filters, 1 per input map
- ► *M*×*N* total filters

#### **The Convolutional Layer: A Closer Look**



#### **Learning Complex Features with CNNs**

Deep CNNs combine simple features into complex patterns

- Early conv layers = edges, textures, ridges
- Later conv layers = eyes, noses, mouths



Image credit: https://devblogs.nvidia.com/parallelforall/accelerate-machine-learning-cudnn-deep-neural-network-library/; H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Unsupervised Learning of Hierarchical Representations with Convolutional Deep Belief Networks", CACM Oct 2011

Part 4
TRANSFORMERS

### **Attention: Local to Global**

- Local pixel relationships make convolutions natural for images
  - Text requires global relationships
  - Attention computes global pairwise relationships



- Efficiency important since global requires O(L<sup>2</sup>) in the input length L
  - Opportunities for domainspecific acceleration
  - Many emerging accelerators for transformers



### **Attention Mechanism**

- The attention mechanism finds pair relationships between all inputs
  - Compute queries, keys, values for each input word
  - Query: "what each word is looking for"
  - Key: "what each word provides"



https://jalammar.github.io/illustrated-transformer/

### **Transformers**

- Combines attention and MLP layers
  - Attention between words
  - MLP within words
- Repeat structure many times in deep network

- Auto-regressive sampling
  - Sample one, send back to network, sample again
  - Based on output probabilities
  - Continue until "end" word



#### **Acknowledgements**

- The lecture slides contain/adapt materials from
  - Dr. Ritchie Zhao (NVIDIA)
  - ECE 5545 by Prof. Mohamed Abdelfattah (Cornell Tech)
  - CS898 by Prof. Ming Li (University of Waterloo)
  - System for AI Education Resource by Microsoft Research