

IEOR 4742 – Deep Learning for OR & FE

Introduction to Feedforward Neural Networks

Ali Hirsa

Industrial Engineering & Operations Research
Columbia University

For educational purposes only, references are not fully cited, some images may be subject to copyright

Hype vs. Reality

- be careful about the hype
- often some most cited papers contain non-reproducible results

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Basics

- basic intuition for ML and NN
- basic structure of a neuron
- how feedforward NNs work
- importance of nonlinearity in tackling complex learning problems

For educational purposes only. References are not cited, some images may be subject to copyright.

Deep Learning (1st definition)

- for decades, scientists have dreamed of making/building intelligent machines w/ brains like human beings (AI)
 - e.g.
 - ✓ self-driving cars
 - ✓ robot housekeepers
 - :
- making these AI machines requires to solve some very computationally complex problems
- to tackle these problems, there is a need to develop a radically different way of programming
- this active field of **Artificial Computer Intelligence** is referred (by some practitioners) to as **Deep Learning**

For educational purposes only. References are not fully cited, some images may be subject to copyright.

Traditional computer programs

traditional computer programs do the following tasks (really) well¹

- (a) performing arithmetic extremely fast
- (b) number crunching
- (c) following list of instructions

For educational purposes only, references not fully cited, some images may be subject to copyright.

¹still very important in trading, market-making, etc

Innovative programs (1 of 4)

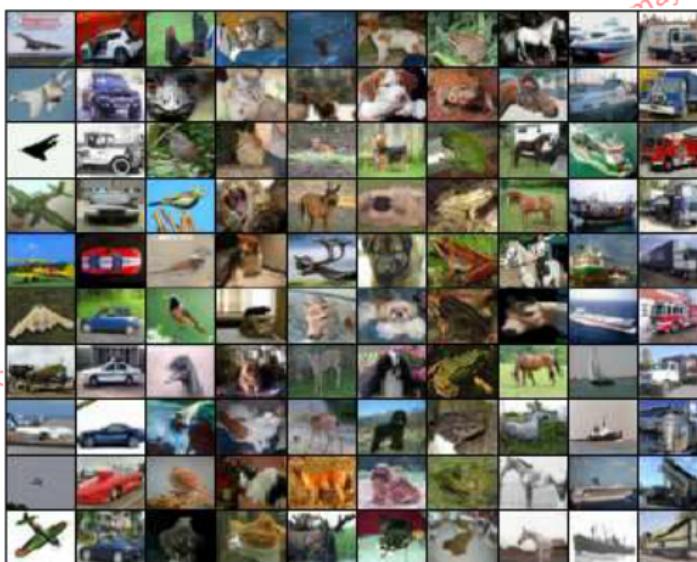
how about writing a program to automatically read yours or someone else's handwriting



Figure: MNIST handwritten digit dataset collected by LeCunn, Cortes, and Burges

Innovative programs (2 of 4)

how about writing a program to automatically recognize an image (an animal, an airplane, etc)



For educational purposes

may be subject to copyright

Figure: CIFAR-10 dataset collected by Krizhevsky, Nair, and Hinton

Innovative programs (3 of 4)

- humans can easily recognize every digit even though they are written in a slightly different way
- how about writing a code/module that achieves this!

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Innovative programs (4 of 4)

- first task is working on MNIST handwritten digit dataset (10 labels) and move to CIFAR-10 image recognition
 - the MNIST database has a training set of 60,000 examples, and a test set of 10,000 examples
 - it is a subset of a larger set available from NIST
 - the digits have been size-normalized and centered in a fixed-size image
 - the CIFAR-10 dataset consists of 60,000 32×32 color images in 10 classes, with 6,000 images per class
 - there are 50,000 training images and 10,000 test images

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Naive Approach (1 of 2)

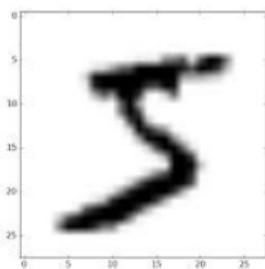
- but before that, we need to set up some rules
- for zero, we might say a closed loop. Is it sufficient? NO



For educational purposes only, references?

- ★ sloppy zero or messy six?
- ★ not a closed loop, how to distinguish?
- ★ some sort of cutoff? for the distance?

Naive Approach (2 of 2)



| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 4 | 9 | 4 | 9 | 4 | 9 | 4 | 9 | 4 | 9 | 4 | 9 |
| 4 | 9 | 7 | 9 | 4 | 9 | 9 | 9 | 4 | 9 | 4 | 9 |
| 9 | 9 | 4 | 9 | 7 | 9 | 4 | 9 | 7 | 9 | 4 | 9 |
| 4 | 9 | 7 | 9 | 4 | 9 | 4 | 9 | 4 | 9 | 4 | 9 |
| 4 | 9 | 7 | 9 | 7 | 9 | 9 | 9 | 7 | 9 | 4 | 9 |
| 7 | 9 | 4 | 9 | 9 | 9 | 4 | 9 | 7 | 9 | 4 | 9 |

For educational purposes only. References are .c

Want to copyrig...

- how about 3 & 5?
- how about 4 & 9?
- we can add more & more rules or **features** but not an easy process
- could be pretty subjective
- AND beginning of our worries!

Industry-Level Tasks (1 of 2)

other classes of problems that fall into this same category:

- object recognition
- speech comprehension
- automated translation/NLP
- pattern recognition

For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Industry-Level Tasks (2 of 2)

signal-to-noise (TR vs. MR)

INTENTIONALLY LEFT BLANK

For educational purposes only, references are not fully cited, some images may be subject to copyright.

uDu vs. UdU

Mechanics of Learning

- we do not exactly know how it is done by our brains
- even if we know, the **program** might be super complicated
- learning by example not by formula.
 - think of a 2-year-old kid who would recognize a dog or a cat by being shown multiple examples of them
 - ★ being corrected if she/he made a wrong guess (if that guess was confirmed by her parents, the learning would be **reinforced**).
 - think of traditional traders, learned simply by simply watching the markets and using charts
 - ★ if markets proved them wrong, modify the thinking to incorporate this **new** information and improve performance

For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Machine Learning vs. Brute-Force Learning

- in ML, instead of teaching a computer a huge list of features/rules to solve the problem, we give it a model with which it can evaluate examples, and a small(er) set of instructions to modify the model when it makes a mistake.

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Short list of Models

Supervised learning (data & **labels**):

- Decision Trees (DTs)
- Classification (Bayes)
- Regression (OLS)
- Logistic Regression
- Support Vector Machine (SVM)
- Ensemble Methods (BAGGing, Random Forest)

Unsupervised Learning (data):

- Clustering
- PCA/SVD
- ICA

For educational purposes only, references not fully cited, some images may be subject to copyright.

Supervised vs. Unsupervised Learning

- supervised learning: given both inputs and labels (outputs)
- typically done in the context of classification, when we want to map input to output labels
- unsupervised learning: learn the inherent structure of our data without using explicitly-provided labels

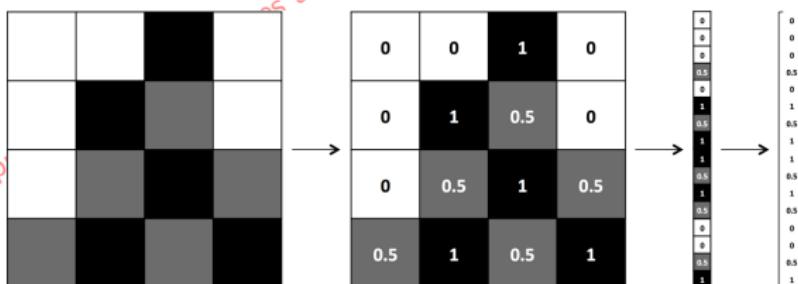
For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Model as a Function

mathematically, a model could be a function:

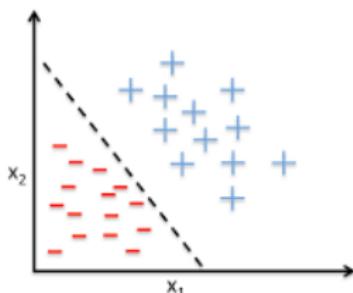
$$h(x, \Theta)$$

In case of x being an image should be converted into a scalar vector



For educational purposes only, some images may be subject to copyright.

Linear decision boundary for classification



Example of a linear decision boundary
for binary classification.

For educational purposes only.

$$h(x, \Theta) = \begin{cases} -1 & \text{if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 < 0 \\ +1 & \text{if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 > 0 \end{cases}$$

Linear Perceptron (1 of 2)

- a model like that is called **Linear Perceptron**

Definition

Perceptron: a computer model or computerized machine devised to represent or simulate the ability of the brain to recognize and discriminate.

- first model for learning with a teacher (i.e., supervised learning) proposed by Frank Rosenblatt 1958

For educational purposes only

Linear Perceptron (2 of 2)

Q: how do we come up w/ optimal value for parameters

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Linear Perceptron (2 of 2)

- Q: how do we come up w/ optimal value for parameters
A: by optimizing a loss function

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Linear Perceptron (2 of 2)

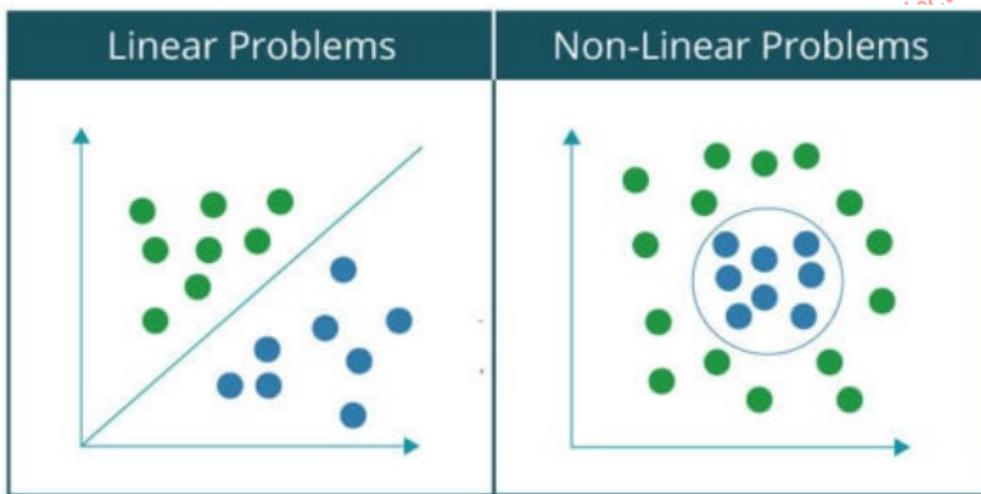
Q: how do we come up w/ optimal value for parameters

A: by optimizing a **loss** function

Goal: to maximize the performance of a machine learning model by iteratively tweaking its parameters until the error is minimized

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Limitation of Linear Perception



For educational purposes only

- Linear perception model is limited
- more complex form is needed

Deep Learning (2nd definition)

to accommodate this complexity:

- researchers have tried to build models that resemble the structures utilized by our brains
- it is essentially this body of research commonly referred to as **Deep Learning**
- **claim:** these algorithms not only far surpass other kinds of ML algos, but also rival the accuracies achieved by humans

For educational purposes only, references are not fully cited, some images may be subject to copyright.

The Neuron

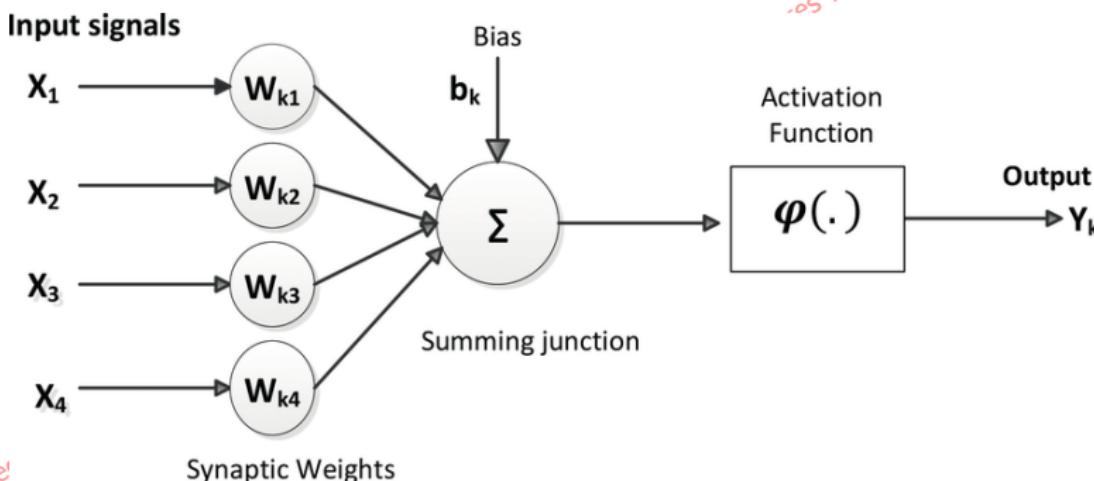
- foundation unit of the human brain is the *neuron*
- Intel claims by 2026 processors will have as many transistors as there are neurons in a brain (100,000,000,000).
- ★ a grain of rice contains over 10,000 neurons, each of which forms an average of 6000 connections w/ other neurons
- this *Massive Biological Network* enables us to experience the world around us.

For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Translation (1 of 2)

McCulloch & Pitts 1943

→ S may be subject to copyright.



For e'

Translation (2 of 2)

translation of this functional understanding of the neurons in our brain into an artificial model that we can represent on our computer is as follows:

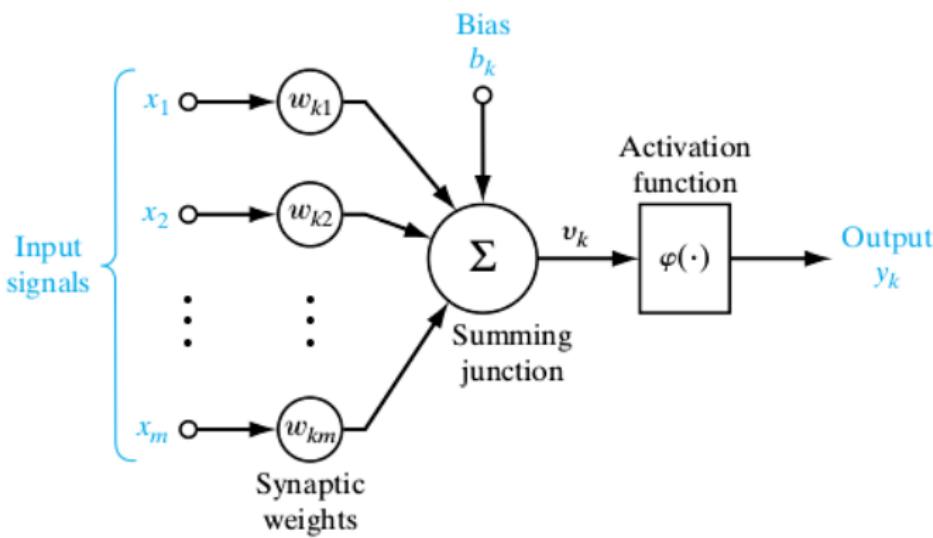
- input x_1, \dots, x_n any non-normalized probability distribution
- weight of the neuron w_1, \dots, w_n
- $y = \varphi(x^\top w + b)$ where b is the biased term

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Linear Perceptrons as Neurons (1 of 2)

$$h(x, \Theta) = \begin{cases} -1 & : w_0 + w_1x_1 + w_2x_2 < 0 \\ 1 & : w_0 + w_1x_1 + w_2x_2 \geq 0 \end{cases}$$

images may be subject to copyright



Linear Perceptrons as Neurons (2 of 2)

- brain is made of more than one neuron

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Linear Perceptrons as Neurons (2 of 2)

- brain is made of more than one neuron
- it is **impossible** for a single neuron to differentiate handwritten digits

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Linear Perceptrons as Neurons (2 of 2)

- brain is made of more than one neuron
- it is **impossible** for a single neuron to differentiate handwritten digits
- neurons in our brain are organized in layers

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Linear Perceptrons as Neurons (2 of 2)

- brain is made of more than one neuron
 - it is **impossible** for a single neuron to differentiate handwritten digits
 - neurons in our brain are organized in layers
- Feedforward Neural Networks
- For educational purposes only, references are not fully cited, some images may be subject to copyright.*

Layers

according to V.B. Mountcastle (1918 - 2015) Neuroscientist at Johns Hopkins (1957):

- Human Cerebral Cortex is made up of six layers
- information flows from one layer to another until sensory input is converted into conceptual understanding
- it is processed by each layer and passed on to the next layer until in the six layer, we conclude whether we are looking at a cat or a dog

For educational purposes only. References are not fully cited, some images may be subject to copyright.

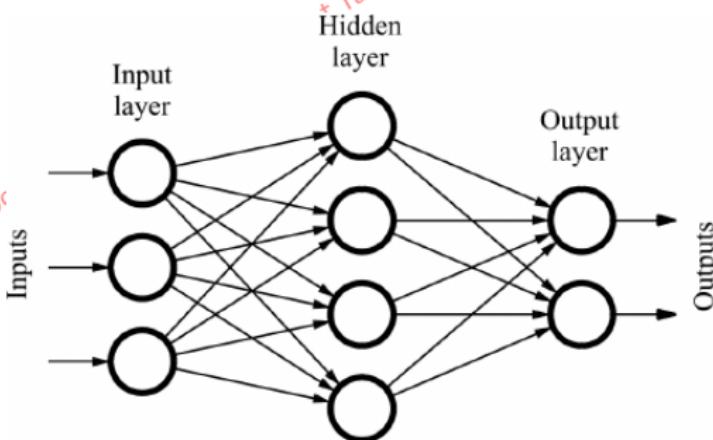
Feedforward Networks (1 of 3)

- a feedforward neural network is an artificial neural network wherein connections between the nodes do not form a cycle
- the feedforward neural network was the first and simplest type of artificial neural network devised
- in this network, the information moves in only one direction, forward, from the input nodes, through the hidden nodes (if any) and to the output nodes
- there are no cycles or loops in the network

For educational purposes only. References are not fully cited. The images may be subject to copyright.

Feedforward Networks (2 of 3)

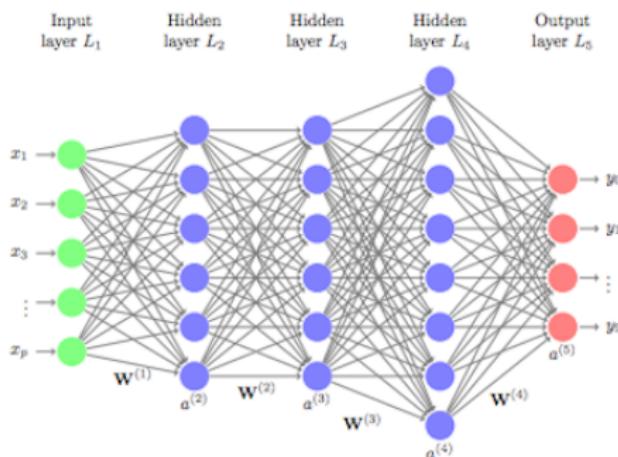
- inputs & outputs are represented as vectors
- connections only travel from lower to higher layer
- no connections in the same layer (for simplicity)
- make it simple to analyze



For educational purpose

Feedforward Networks (3 of 3)

- weight connecting i^{th} neuron in the k^{th} layer with j^{th} neuron in the $(k+1)^{th}$ layer
- these weights constitute the parameter set Θ (vector)
- our ability to solve problems with neural networks depends on finding the optimal values



Hidden Layers

- layers of neurons that are between the input (first layer) and the output (last layer) are called **hidden layers**
- each hidden layer could have different numbers of neurons
- not required that every neuron has its output connected to the inputs of all neurons in the next layer
- how to do that, it is “an art”
- most action is happening when the neural net tries to solve problems
- taking a look at the activities of hidden layers can say a lot about features
 - the network has automatically learned to extract from data
 - choose between fewer layers with more neurons per layer vs. deep layers with fewer numbers of neurons per layer
- usually no more than 5 to 7 layers, overfitting starts to occur

For educational purposes only, references are not fully cited, some images may be subject to copyright.

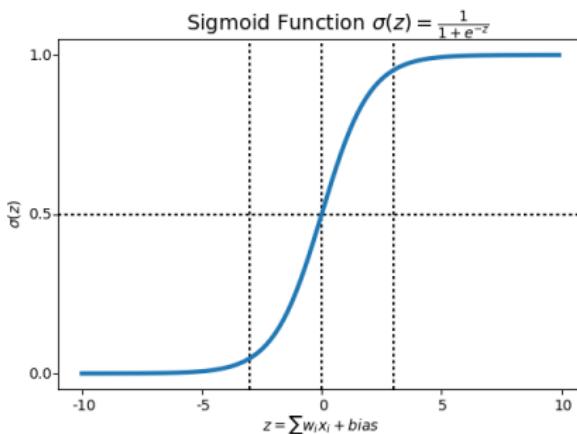
Activation Functions

- linear neurons are easy to compute with, but have serious limitations
- can be shown that any feedforward neural network consisting of only linear neurons can be expressed as a network with no hidden layers (**exercise**)

For educational purposes only references are not fully cited. Some images may be subject to copyright.

The logistic sigmoid function

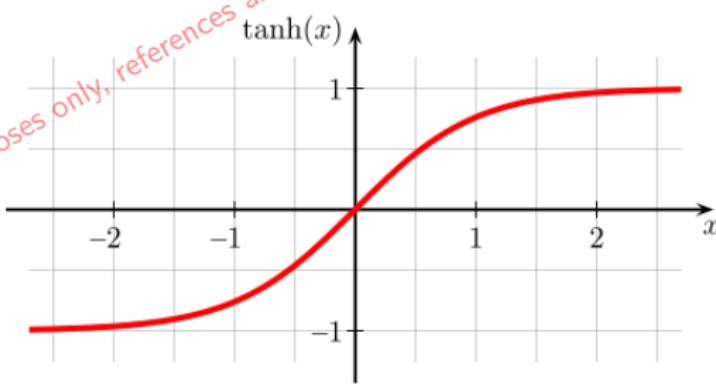
$$\begin{aligned}f(z) &= \frac{e^z}{1 + e^z} \\&= \frac{1}{1 + e^{-z}}\end{aligned}$$



The tanh function

$$\begin{aligned}f(z) &= \tanh(z) \\&= \frac{e^z - e^{-z}}{e^z + e^{-z}} \\&= 2\sigma(2z) - 1\end{aligned}$$

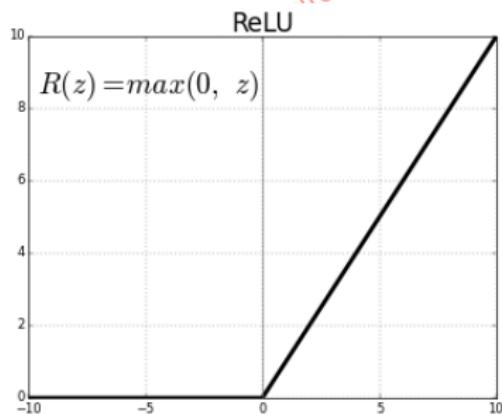
tanh is a re-scaled logistic sigmoid function



For educational purposes only, references are not fully cited, some images may be subject to copyright.

The restricted linear unit (ReLU) function

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

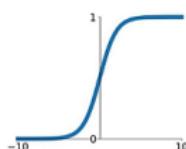


Collection of deep learning activation functions (1 of 2)

may be subject to copyright

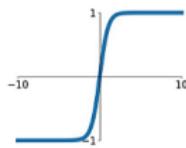
Sigmoid

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



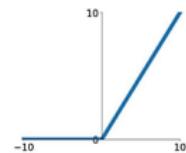
tanh

$$\tanh(x)$$



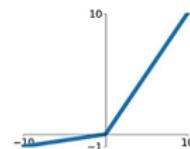
ReLU

$$\max(0, x)$$



Fo'

Leaky ReLU

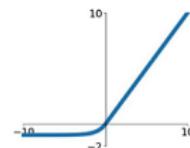
$$\max(0.1x, x)$$


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

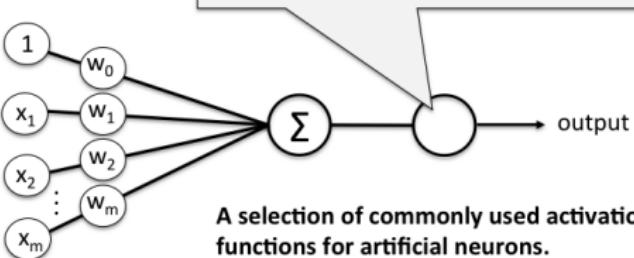
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Collection of deep learning activation functions (2 of 2)

-copyrig

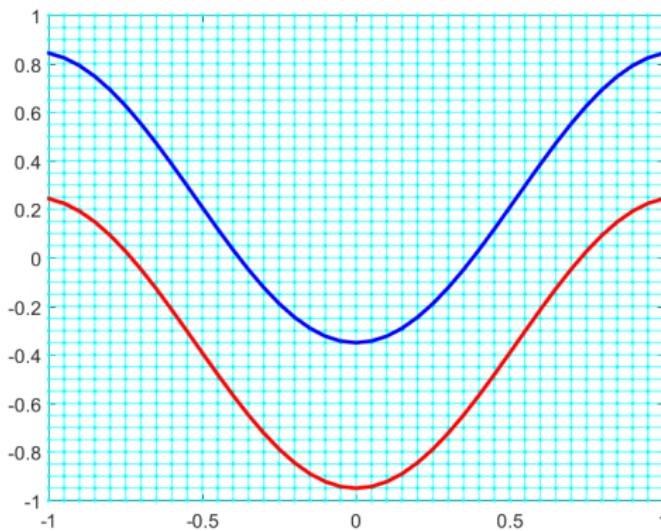
| | | |
|---|------------------------------|--|
|  | Unit step | $g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{otherwise.} \end{cases}$ |
|  | Linear | $g(z) = z$ |
|  | Logistic (sigmoid) | $g(z) = 1 / (1 + \exp(-z))$ |
|  | Hyperbolic tangent (sigmoid) | $g(z) = \frac{\exp(2z) - 1}{\exp(2z) + 1}$ |
| ... | | |



For edv

Linear Classification Example (1 of 3)

- very simple dataset, two curves on a plane



For educational p'

Linear Classification Example (2 of 3)

- to classify points as belonging to one or the other curve
- obvious way to visualize the behavior of any classification algorithm is to simply look at how it classifies every possible data point
- first trial is to simply separate two classes of data by dividing them with a line

For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Linear Classification Example (3 of 3)

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Transforming the data

- as expected, the first trial failed to do the job
- how about a more complicated curve than a line?
- how about transforming the data, and creating a new representation?
- AND once we get to the final representation, network would just draw a line (in higher dimensions, a hyperplane) through the transformed data

For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Topology (1 of 3)

- Topology is the mathematical study of the properties that are preserved through deformation, twisting, and stretching of objects
- tearing, however, is not allowed
- for example, a circle is topologically equivalent to an ellipse (into which it can be deformed by stretching) and a sphere is equivalent to an ellipsoid, and so on

For educational purposes only. References are not fully cited. Some images may be subject to copyright.

Topology (2 of 3)

- definition of topology leads to the following mathematical joke
(Renteln and Dundes 2005)

Q: who is a topologist?

A: someone who cannot distinguish between a doughnut and a coffee cup

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Topology (3 of 3)

- an n -dimensional topological space is a space (not necessarily Euclidean) with certain properties of connectedness and compactness

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Naive approach for visualization purposes

for every grid point in the manifold

Algorithm 1: pseudo-code for transforming the manifold

```
1 for  $w_{11} = -3, \dots, 3$  do
2   for  $w_{12} = -3, \dots, 3$  do
3     for  $w_{21} = -3, \dots, 3$  do
4       for  $w_{22} = -3, \dots, 3$  do
5         for  $b_1 = -1, \dots, 1$  do
6           for  $b_2 = -1, \dots, 1$  do
7              $\hat{x} = \tanh(w_{11}x + w_{21}y + b_1)$ 
8              $\hat{y} = \tanh(w_{12}x + w_{22}y + b_2)$ 
9           end
10        end
11      end
12    end
13  end
14 end
```

For educational purposes only. References are not fully cited, some images may be subject to copyright.

Transformed Manifold

For educational purposes only, references are not fully cited, some images may be subject to copyright.

Loss function

note we seek \mathbf{w} such that

$$\mathbf{w}^\top \mathbf{x} \geq 0 \text{ when } t = +1$$

$$\mathbf{w}^\top \mathbf{x} < 0 \text{ when } t = -1$$

or equivalently

$$\mathbf{w}^\top \mathbf{x}_j t_j \geq 0 \quad \forall j$$

thus we seek to minimize

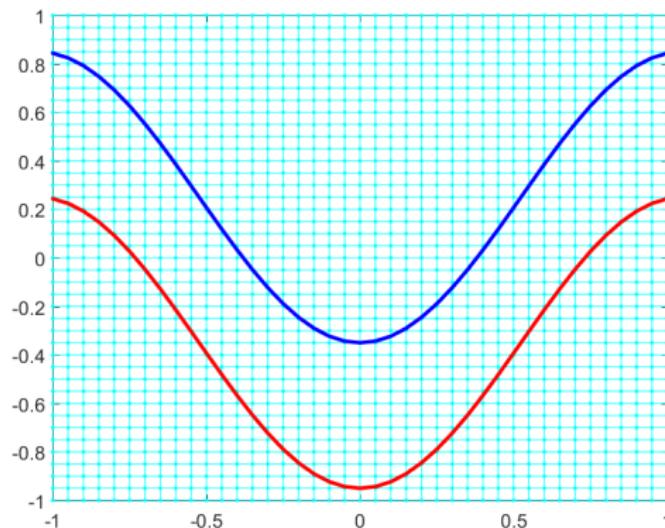
$$-\sum_{j \in \mathcal{A}} \mathbf{w}^\top \mathbf{x}_j t_j \geq 0 \quad \forall j$$

where \mathcal{A} is the set of mis-classified inputs

Linear Classification Example (1 of 2)

Subject to copyright

For educational p'

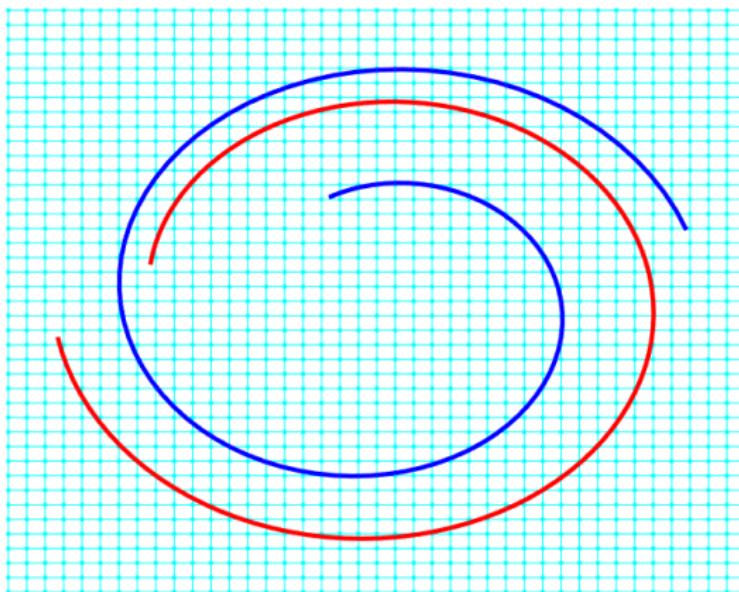


Linear Classification Example (2 of 2)

For educational purposes only, references are not fully cited, some images may be subject to copyright.

A more complicated example

+o copyrig



For educati