

# Deep Feedforward Networks: Overview

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Overview
- 1.Example: Learning XOR
- 2.Gradient-Based Learning
- 3.Hidden Units
- 4.Architecture Design
- 5.Backpropagation and Other Differentiation
- 6.Historical Notes

# Goal of a feedforward network

- Feedforward Nets are quintessential deep learning models
- Deep Feedforward Networks are also called as
  - Feedforward neural networks or
  - Multilayer Perceptrons (MLPs)
- Their Goal is to approximate some function  $f^*$ 
  - E.g., classifier  $y = f^*(x)$  maps input  $x$  to category  $y$
  - Feedforward Network defines a mapping

$$y = f^*(x; \theta)$$

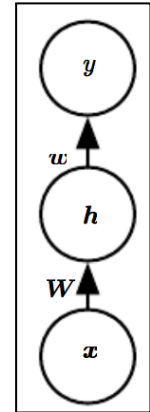
- and learns the values of the parameters  $\theta$  that result in the best function approximation

# Flow of Information

- Models are called *Feedforward* because:

- To evaluate  $f(\mathbf{x})$ : information flows one-way from  $\mathbf{x}$  through computations defining  $f$  to output  $y$

$$y = f(\mathbf{x})$$



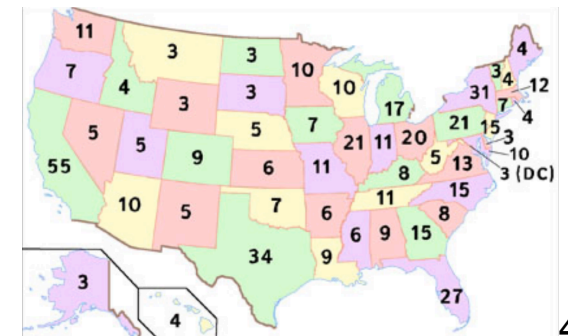
- There are no feedback connections

- No outputs of model are fed back into itself

- US Presidential Election

- Input:  $\mathbf{x} = \{x_1, \dots, x_{50}\}$ 
  - are votes cast for a candidate
- $h$  is electoral college
  - Each state has fixed no of electors
- Output:  $y$ 
  - votes of electoral college for candidate

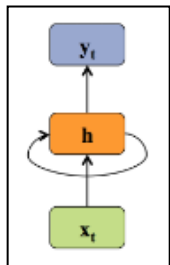
$h$  is defined for each state as shown in map



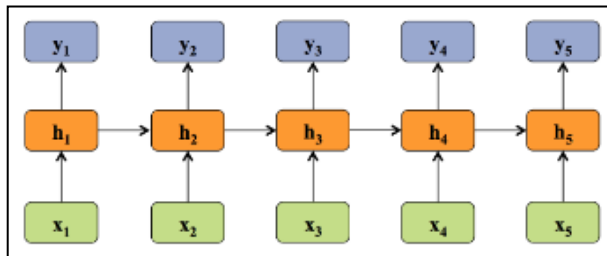
# Feedforward vs. Recurrent

- When feedforward neural networks are extended to include feedback connections they are called *Recurrent Neural Networks (RNNs)*

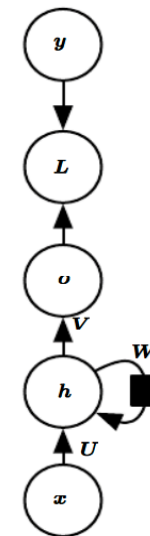
RNN



Unrolled RNN



RNN with  
learning  
component



# Importance of Feedforward Networks

- They are extremely important to ML practice
- Form basis for many commercial applications
  1. CNNs are a special kind of feedforward networks
    - They are used for recognizing objects from photos
  2. They are a conceptual stepping stones to RNNs
    - RNNs power many NLP applications

# Feedforward Neural Network Structures

- They are called networks because they are composed of many different functions
- Model is associated with a directed acyclic graph describing how functions composed
  - E.g., functions  $f^{(1)}$ ,  $f^{(2)}$ ,  $f^{(3)}$  connected in a chain to form  $f(\mathbf{x}) = f^{(3)}[f^{(2)}[f^{(1)}(\mathbf{x})]]$ 
    - $f^{(1)}$  is called the first layer of network (which is a vector)
    - $f^{(2)}$  is called the second layer, etc
- These chain structures are the most commonly used structures of neural networks

# Definition of Depth

- Overall length of the chain is the *depth* of the model
  - Ex: the composite function  $f(x) = f^{(3)} [ f^{(2)} [ f^{(1)}(x) ] ]$  has depth of 3
- The name *deep learning* arises from this terminology
- Final layer of a feedforward network, ex  $f^{(3)}$ , is called the *output layer*



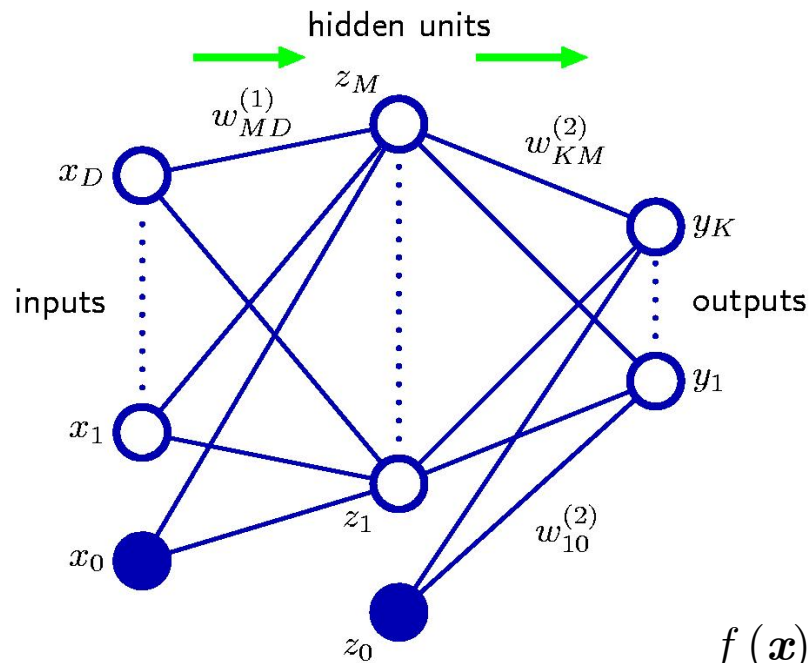
# Training the Network

- In network training we drive  $f(\mathbf{x})$  to match  $f^*(\mathbf{x})$
- Training data provides us with noisy, approximate examples of  $f^*(\mathbf{x})$  evaluated at different training points
- Each example accompanied by label  $y \approx f^*(\mathbf{x})$
- Training examples specify directly what the output layer must do at each point  $\mathbf{x}$ 
  - It must produce a value that is close to  $y$

# Definition of Hidden Layer

- Behavior of other layers is not directly specified by the data
- Learning algorithm must decide how to use those layers to produce value that is close to  $y$
- Training data does not say what individual layers should do
- Since the desired output for these layers is not shown, they are called *hidden layers*

# A net with depth 2: one hidden layer



$K$  outputs  $y_1, \dots, y_K$  for a given input  $\mathbf{x}$   
 Hidden layer consists of  $M$  units

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{j=1}^M w_{kj}^{(2)} h \left( \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

$$f(\mathbf{x}) = f^{(2)} [f^{(1)}(\mathbf{x})]$$

$f^{(1)}$  is a vector of  $M$  dimensions and  
 $f^{(2)}$  is a vector of  $K$  dimensions

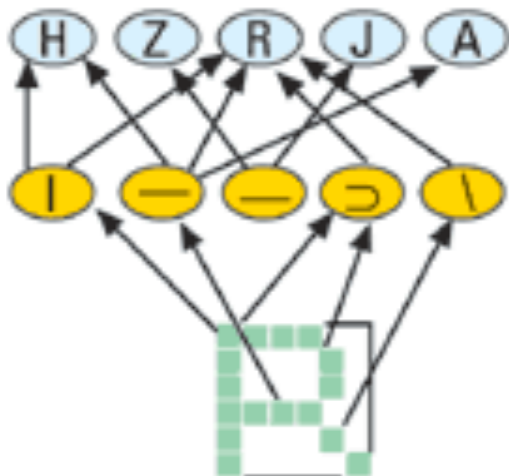
$$\begin{aligned} f_m^{(1)} &= z_m = h(\mathbf{x}^T \mathbf{w}^{(1)}), \quad m=1, \dots, M \\ f_k^{(2)} &= \sigma(\mathbf{z}^T \mathbf{w}^{(2)}), \quad k=1, \dots, K \end{aligned}$$

# Feedforward net with depth 2

- Recognition of printed characters (OCR)

$$f(\mathbf{x}) = f^{(2)}[f^{(1)}(\mathbf{x})]$$

- Hidden layer  $f^{(1)}$  compares raw pixel inputs to component patterns



# Width of Model

- Each hidden layer is typically vector-valued
- Dimensionality of hidden layer vector is *width* of the model

# Units of a model

- Each element of vector viewed as a neuron
  - Instead of thinking of it as a vector-vector function, they are regarded as units in parallel
- Each unit receives inputs from many other units and computes its own activation value

# Depth versus Width

- Going deeper makes network more expressive
  - It can capture variations of the data better.
  - Yields expressiveness more efficiently than width
- Tradeoff for more expressiveness is increased tendency to overfit
  - You will need more data or additional regularization
    - network should be as deep as training data allows.
  - But you can only determine a suitable depth by experiment.
    - Also computation increases with no. of layers.

# Very Deep CNNs

## CNNs with depth 11 to 19

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Depth increases from left (A) to right (E)  
as more layers are added  
(the added layers are shown in bold)

Convolutional layer parameters denoted  
“conv (receptive field size) –(no. of channels)”

ReLU activation not shown for brevity

Table 2: Number of parameters (in millions).

Network	A, A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144



# Why are they neural networks?

- These networks are loosely inspired by neuroscience
- Each unit resembles a neuron
  - Receives input from many other units
  - Computes its own activation value
- Choice of functions  $f^{(i)}(\mathbf{x})$ :
  - Loosely guided by neuroscientific observations about biological neurons
    - Modern neural networks are guided by many mathematical and engineering disciplines
    - Not perfectly model the brain

# Function Approximation is goal

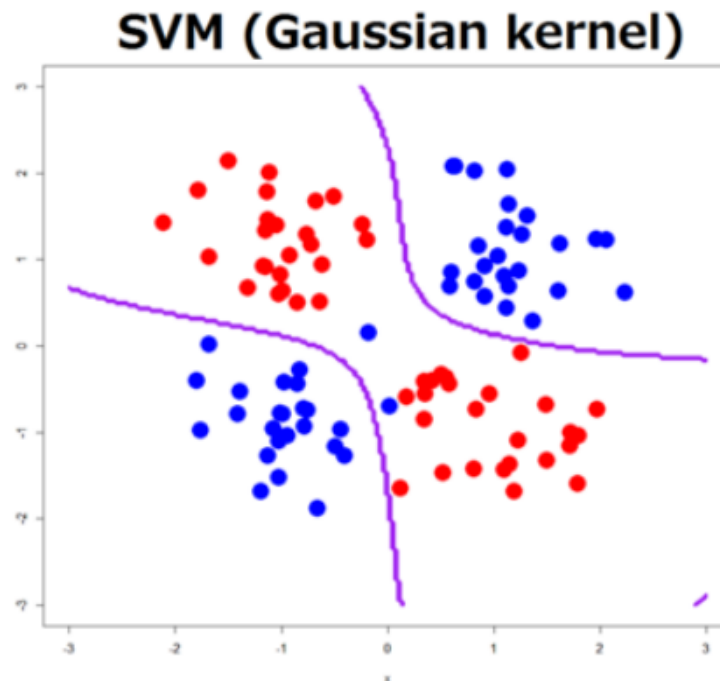
- Think of feedforward networks as function approximation machines
  - Designed to achieve statistical generalization
- Occasionally draw insights from what we know about the brain
  - Rather than as models of brain function

# Understanding Feedforward Nets

- Begin with linear networks and understand their limitations
- Linear models such as logistic regression and linear regression can be fit reliably and efficiently using either
  - Closed-form solution
  - Convex optimization
- Limitation

# Extending Linear Models

- To represent non-linear functions of  $x$ 
  - apply linear model to transformed input  $\phi(x)$ 
    - where  $\phi$  is non-linear
  - Equivalently kernel trick of SVM obtains nonlinearity



# SVM Kernel trick

- Many ML algos can be rewritten as dot products between examples:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \text{ written as } b + \sum_i \alpha_i \mathbf{x}^T \mathbf{x}^{(i)}$$

where  $\mathbf{x}^{(i)}$  is a training example and  $\boldsymbol{\alpha}$  is a vector of coeffs

- This allows us to replace  $\mathbf{x}$  with a feature function  $\phi(\mathbf{x})$  and dot product with function  $k(\mathbf{x}, \mathbf{x}^{(i)}) = \phi(\mathbf{x}) \bullet \phi(\mathbf{x}^{(i)})$  called a kernel
  - The  $\bullet$  operator represents an inner product analogous to  $\phi(\mathbf{x})^T \phi(\mathbf{x}^{(i)})$
  - For some feature spaces we may not literally use an inner product
    - In continuous spaces an inner product based on integration
- Gaussian kernel
  - Consider  $k(\mathbf{u}, \mathbf{v}) = \exp(-\|\mathbf{u} - \mathbf{v}\|^2 / 2\sigma^2)$ 
    - By expanding the square  $\|\mathbf{u} - \mathbf{v}\|^2 = \mathbf{u}^T \mathbf{u} + \mathbf{v}^T \mathbf{v} - 2\mathbf{u}^T \mathbf{v}$
    - we get  $k(\mathbf{u}, \mathbf{v}) = \exp(-\mathbf{u}^T \mathbf{u} / 2\sigma^2) \exp(-\mathbf{u}^T \mathbf{v} / \sigma^2) \exp(-\mathbf{v}^T \mathbf{v} / 2\sigma^2)$
  - Validity follows from kernel construction rules

# SVM Prediction

- Use linear regression on Lagrangian for determining the weights  $\alpha_i$
- We can make predictions using
  - $f(\mathbf{x}) = b + \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$
  - Function is nonlinear wrt  $\mathbf{x}$  but relationship between  $\phi(\mathbf{x})$  and  $f(\mathbf{x})$  is linear
  - Also the relationship between  $\alpha$  and  $f(\mathbf{x})$  is linear
  - We can think of  $\phi$  as providing a set of features
    - describing  $\mathbf{x}$  or providing a new representation for  $\mathbf{x}$

# Disadvantages of Kernel Methods

- Cost of decision function evaluation: linear in  $m$ 
  - Because the  $i^{\text{th}}$  example contributes a term  $\alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$  to the decision function
  - Can mitigate this by learning an  $\alpha$  with mostly zeros
    - Classification requires evaluating the kernel function only for training examples that have a nonzero  $\alpha_i$
    - These are known as *support vectors*
- Cost of training: high with large data sets
- Generic kernels struggle to generalize well
  - Neural net outperformed RBF-SVM on MNIST
- Also, how to choose the mapping  $\phi$ ?

# Options for choosing mapping $\phi$

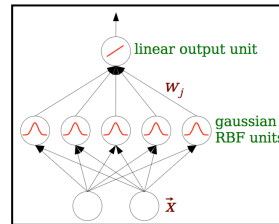
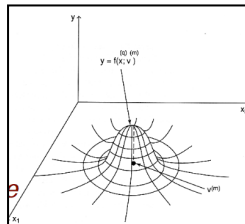
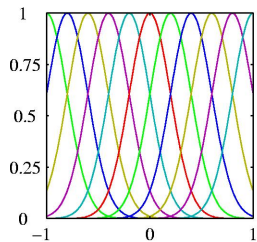
1. Generic feature function  $\phi(x)$ 
  - Radial basis function
2. Manually engineer  $\phi$ 
  - Feature engineering
3. Principle of Deep Learning: Learn  $\phi$



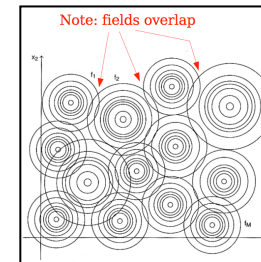
# Option 1 to choose the mapping $\phi$

- Generic feature function  $\phi(x)$ 
  - Infinite-dimensional  $\phi$  that is implicitly used by kernel machines based on RBF

- RBF:  $N(x; x^{(i)}, \sigma^2 I)$  centered at  $x^{(i)}$



$x^{(i)}$ : From  $k$ -means clustering



$\sigma$  = mean distance between each unit  $j$  and its closest neighbor

- If  $\phi(x)$  is of high enough dimension we can have enough capacity to fit the training set
  - Generalization to test set remains poor
  - Generic feature mappings are based on smoothness
    - Do not include prior information to solve advanced problems

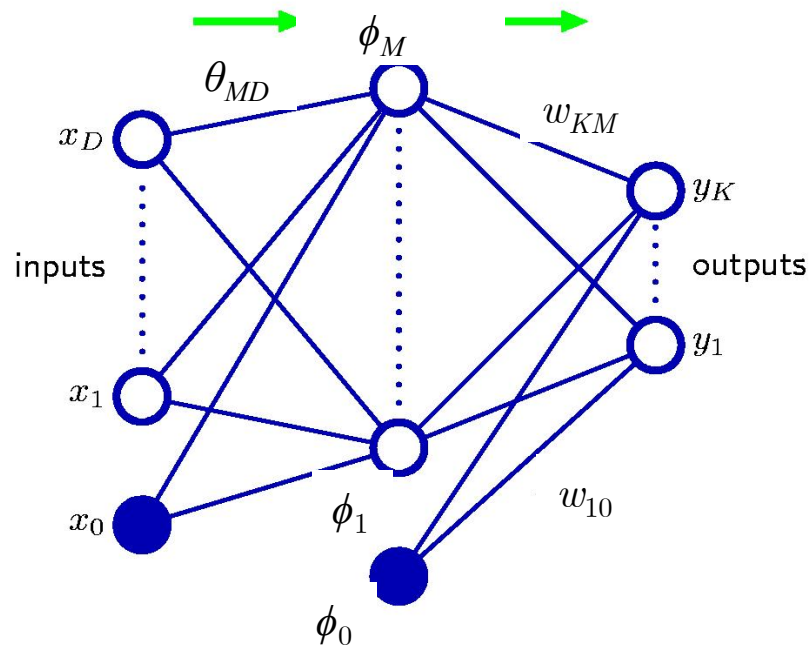
# Option 2 to choose the mapping $\phi$

- Manually engineer  $\phi$
- This was the dominant approach until arrival of deep learning
- Requires decades of effort
  - e.g., speech recognition, computer vision
- Little transfer between domains

# Option 3 to choose the mapping $\phi$

- Strategy of Deep learning: Learn  $\phi$
- Model is  $y = f(x; \theta, w) = \phi(x; \theta)^T w$ 
  - $\theta$  used to learn  $\phi$  from broad class of functions
  - Parameters  $w$  map from  $\phi(x)$  to output
  - Defines FFN where  $\phi$  define a hidden layer
- Unlike other two (basis functions, manual engineering), this approach gives-up on convexity of training
  - But its benefits outweigh harms

# Extend Linear Methods to Learn $\phi$



$K$  outputs  $y_1, \dots, y_K$  for a given input  $\mathbf{x}$   
 Hidden layer consists of  $M$  units

$$y_k(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \sum_{j=1}^M w_{kj} \phi_j \left( \sum_{i=1}^D \theta_{ji} x_i + \theta_{j0} \right) + w_{k0}$$

$$y_k = f_k(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \boldsymbol{\phi}(\mathbf{x}; \boldsymbol{\theta})^T \mathbf{w}$$

Can be viewed as a generalization of linear models

- Nonlinear function  $f_k$  with  $M+1$  parameters  $\mathbf{w}_k = (w_{k0}, \dots, w_{kM})$  with
- $M$  basis functions,  $\phi_j$   $j=1, \dots, M$  each with  $D$  parameters  $\boldsymbol{\theta}_j = (\theta_{j1}, \dots, \theta_{jD})$
- Both  $\mathbf{w}_k$  and  $\boldsymbol{\theta}_j$  are learnt from data

# Approaches to Learning $\phi$

- Parameterize the basis functions as  $\phi(x;\theta)$ 
  - Use optimization to find  $\theta$  that corresponds to a good representation
- Approach can capture benefit of first approach (fixed basis functions) by being highly generic
  - By using a broad family for  $\phi(x;\theta)$
- Can also capture benefits of second approach
  - Human practitioners design families of  $\phi(x;\theta)$  that will perform well
  - Need only find right function family rather than precise right function

# Importance of Learning $\phi$

- Learning  $\phi$  is discussed beyond this first introduction to feed-forward networks
  - It is a recurring theme throughout deep learning applicable to all kinds of models
- Feedforward networks are application of this principle to learning deterministic mappings from  $x$  to  $y$  without feedback
- Applicable to
  - learning stochastic mappings
  - functions with feedback
  - learning probability distributions over a single vector

# Plan of Discussion: Feedforward Networks

1. A simple example: learning XOR
2. Design decisions for a feedforward network
  - Many are same as for designing a linear model
    - Basics of gradient descent
      - Choosing the optimizer, Cost function, Form of output units
  - Some are unique
    - Concept of hidden layer
      - Makes it necessary to have activation functions
    - Architecture of network
      - How many layers , How are they connected to each other, How many units in each later
    - Learning requires gradients of complicated functions
      - Backpropagation and modern generalizations

# 1. Ex: XOR problem

- XOR: an operation on binary variables  $x_1$  and  $x_2$ 
  - When exactly one value equals 1 it returns 1 otherwise it returns 0
  - Target function is  $y=f^*(\mathbf{x})$  that we want to learn
    - Our model is  $y=f([x_1, x_2]; \boldsymbol{\theta})$  which we learn, i.e., adapt parameters  $\boldsymbol{\theta}$  to make it similar to  $f^*$
- Not concerned with statistical generalization
  - Perform correctly on four training points:
    - $X=\{[0,0]^T, [0,1]^T, [1,0]^T, [1,1]^T\}$
  - Challenge is to fit the training set
    - We want  $f([0,0]^T; \boldsymbol{\theta}) = f([1,1]^T; \boldsymbol{\theta}) = 0$
    - $f([0,1]^T; \boldsymbol{\theta}) = f([1,0]^T; \boldsymbol{\theta}) = 1$



# ML for XOR: linear model doesn't fit

- Treat it as regression with MSE loss function

$$J(\theta) = \frac{1}{4} \sum_{\mathbf{x} \in X} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2 = \frac{1}{4} \sum_{n=1}^4 (f^*(\mathbf{x}_n) - f(\mathbf{x}_n; \theta))^2$$

- Usually not used for binary data
- But math is simple

Alternative is Cross-entropy  $J(\theta)$

$$\begin{aligned} J(\theta) &= -\ln p(\mathbf{t} | \theta) \\ &= -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \\ y_n &= \sigma(\theta^T \mathbf{x}_n) \end{aligned}$$

- We must choose the form of the model
- Consider a linear model with  $\theta = \{\mathbf{w}, b\}$  where

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{x}^T \mathbf{w} + b$$

- Minimize  $J(\theta) = \frac{1}{4} \sum_{n=1}^4 (t_n - \mathbf{x}_n^T \mathbf{w} - b)^2$  to get closed-form solution

- Differentiate wrt  $\mathbf{w}$  and  $b$  to obtain  $\mathbf{w} = \mathbf{0}$  and  $b = 1/2$

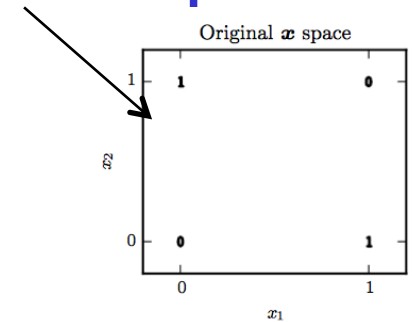
- Then the linear model  $f(\mathbf{x}; \mathbf{w}, b) = 1/2$  simply outputs 0.5 everywhere

- Why does this happen?

# Linear model cannot solve XOR

- Bold numbers are values system must output

- When  $x_1=0$ , output has to increase with  $x_2$
- When  $x_1=1$ , output has to decrease with  $x_2$



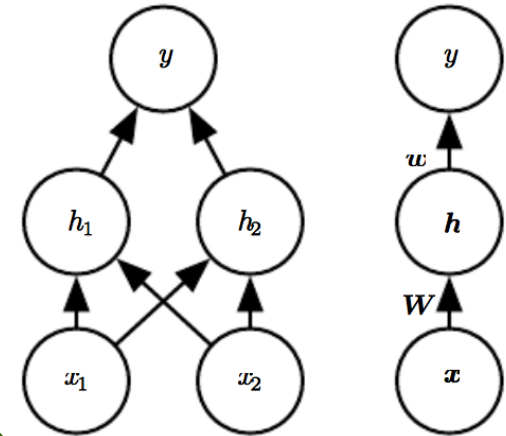
- Linear model  $f(\mathbf{x}; \mathbf{w}, b) = x_1 w_1 + x_2 w_2 + b$  has to assign a single weight to  $x_2$ , so it cannot solve this problem

- A better solution:

- use a model to learn a different representation
  - in which a linear model is able to represent the solution
- We use a simple feedforward network
  - one hidden layer containing two hidden units

# Feedforward Network for XOR

- Introduce a simple feedforward network
  - with one hidden layer containing two units
- Same network drawn in two different styles
  - Matrix  $W$  describes mapping from  $x$  to  $h$
  - Vector  $w$  describes mapping from  $h$  to  $y$
  - Intercept parameters  $b$  are omitted



# Functions computed by Network

- Layer 1 (hidden layer): vector of hidden units  $h$  computed by function  $f^{(1)}(\mathbf{x}; W, \mathbf{c})$ 
  - $\mathbf{c}$  are bias variables

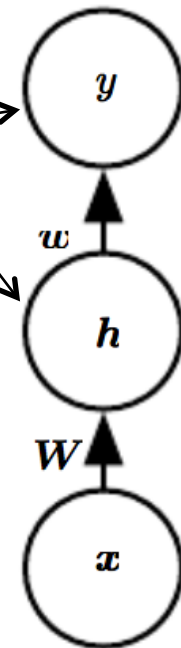
- Layer 2 (output layer) computes

$$f^{(2)}(\mathbf{h}; \mathbf{w}, b)$$

- $\mathbf{w}$  are linear regression weights
- Output is linear regression applied to  $h$  rather than to  $\mathbf{x}$

- Complete model is

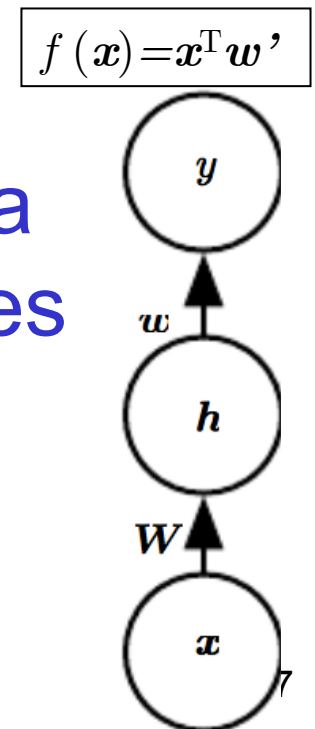
$$f(\mathbf{x}; W, \mathbf{c}, \mathbf{w}, b) = f^{(2)}(f^{(1)}(\mathbf{x}))$$



# Linear vs Nonlinear functions

- If we choose both  $f^{(1)}$  and  $f^{(2)}$  to be linear, the total function will still be linear  $f(x) = x^T w'$ 
  - Suppose  $f^{(1)}(x) = W^T x$  and  $f^{(2)}(h) = h^T w$
  - Then we could represent this function as
- Since linear is insufficient, we must use a nonlinear function to describe the features
  - We use the strategy of neural networks
  - by using a nonlinear activation function

$$h = g(W^T x + c)$$



# Activation Function

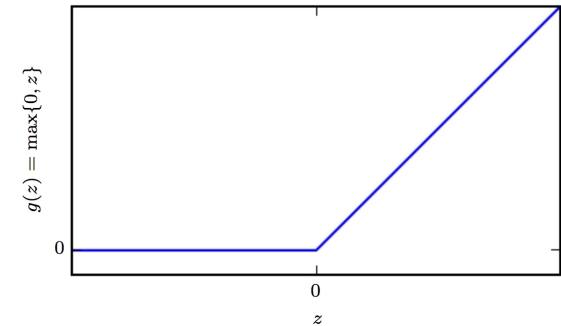
- In linear regression we used a vector of weights  $w$  and scalar bias  $b$ 

$f(x; w, b) = x^T w + b$

  - to describe an affine transformation from an input vector to an output scalar
- Now we describe an affine transformation from a vector  $x$  to a vector  $h$ , so an entire vector of bias parameters is needed
- Activation function  $g$  is typically chosen to be applied element-wise  $h_i = g(x^T W_{:,i} + c_i)$

# Default Activation Function

- **Activation:**  $g(z) = \max\{0, z\}$ 
  - Applying this to the output of a linear transformation yields a nonlinear transformation
  - However function remains close to linear
    - Piecewise linear with two pieces
    - Therefore preserve properties that make linear models easy to optimize with gradient-based methods
    - Preserve many properties that make linear models generalize well



## A principle of CS:

Build complicated systems from minimal components.

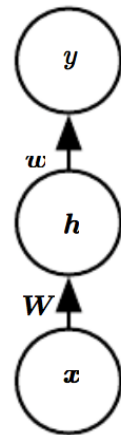
A Turing Machine  
Memory needs only 0 and 1 states.

We can build Universal Function approximator from ReLUs

# Specifying the Network using ReLU

- Activation:  $g(z) = \max\{0, z\}$
- We can now specify the complete network as

$$f(\mathbf{x}; W, \mathbf{c}, \mathbf{w}, \mathbf{b}) = f^{(2)}(f^{(1)}(\mathbf{x})) = \mathbf{w}^T \max\{0, W^T \mathbf{x} + \mathbf{c}\} + \mathbf{b}$$

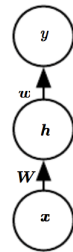




# We can now specify XOR Solution

- Let  $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$
- Now walk through how model processes a batch of inputs

$$f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$$



- Design matrix  $X$  of all four points:

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

- First step is  $XW$ :

$$XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

- Adding  $c$ : In this space all points lie along a line with slope 1. Cannot be implemented by a linear model

$$XW + c = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

- Compute  $h$  Using ReLU

Has changed relationship among examples. They no longer lie on a single line. A linear model suffices

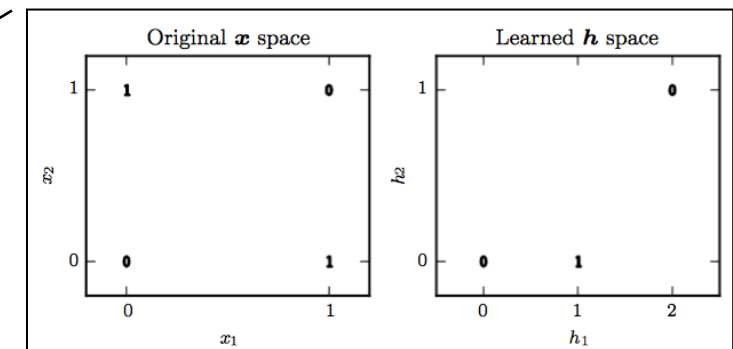
$$\max\{0, XW + c\} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

- Finish by multiplying by  $w$ :

$$f(x) = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

- Network has obtained

correct answer for all 4 examples



# Learned representation for XOR

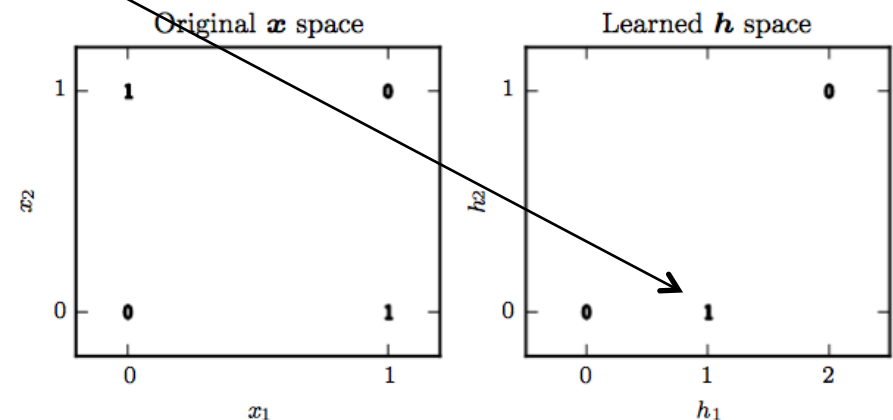
- Two points that must have output 1 have been collapsed into one

- Points  $\mathbf{x}=[0,1]^T$  and  $\mathbf{x}=[1,0]^T$  have been mapped into  $\mathbf{h}=[0,1]^T$

- Described in linear model
  - For fixed  $h_2$ , output increases in  $h_1$

When  $x_1=0$ , output has to increase with  $x_2$

When  $x_1=1$ , output has to decrease with  $x_2$



When  $h_1=0$ , output is constant 0 with  $h_2$

When  $h_1=1$ , output is constant 1 with  $h_2$

When  $h_1=2$ , output is constant 0 with  $h_2$

# About the XOR example

- We simply specified the solution
  - Then showed that it achieves zero error
- In real situations there might be billions of parameters and billions of training examples
  - So one cannot simply guess the solution
- Instead gradient descent optimization can find parameters that produce very little error
  - The solution described is at the global minimum
    - Gradient descent could converge to this solution
    - Convergence depends on initial values
    - Would not always find easily understood integer solutions