# CS489/698: Intro to ML

Lecture 03: Multi-layer Perceptron
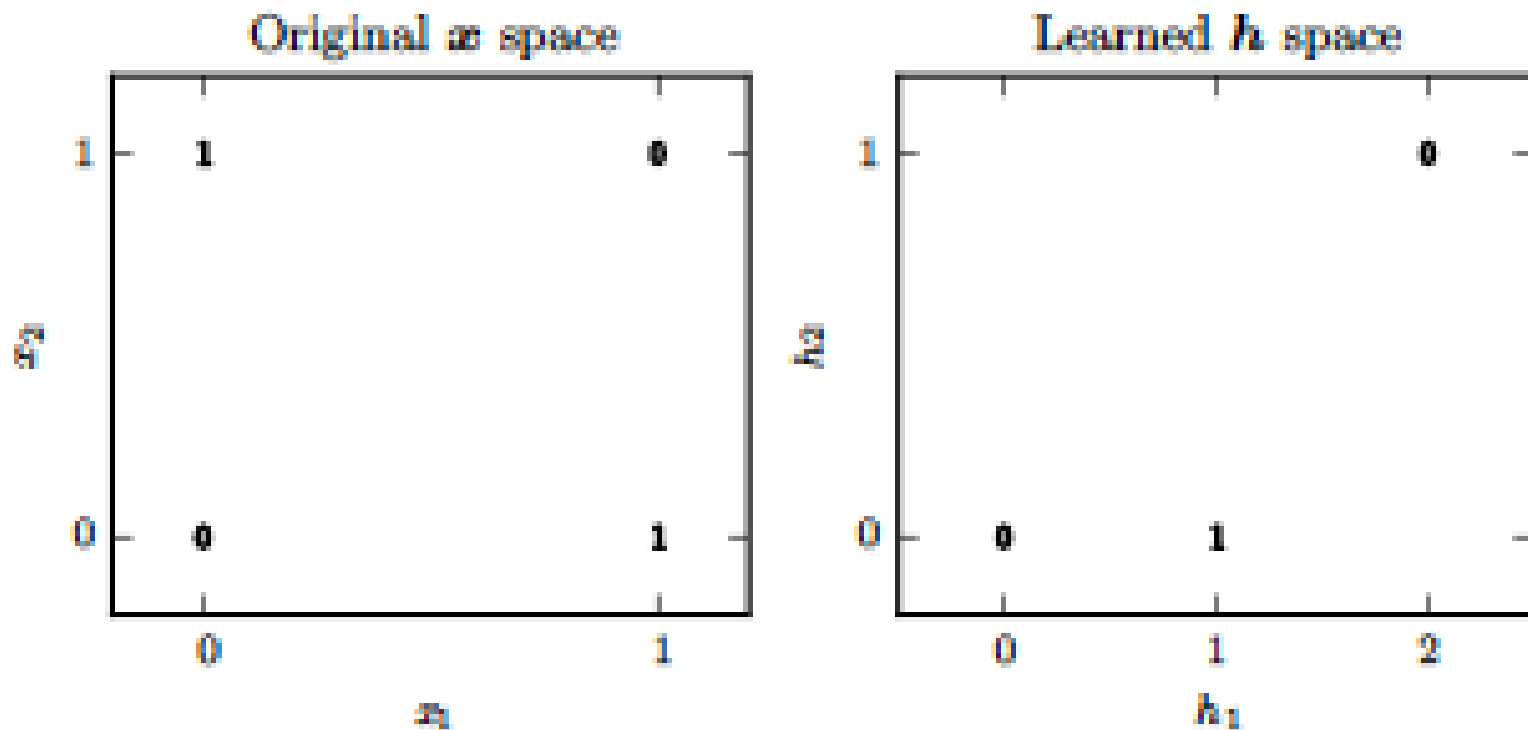
Yao-Liang Yu

# Outline

- **Failure of Perceptron**

- Neural Network

- Backpropagation

- Universal Approximator

9/19/17                                    Yao-Liang Yu

# Outline

- **Failure of Perceptron**

- Neural Network

- Backpropagation

- Universal Approximator

9/19/17    Yao-Liang Yu

**UNIVERSITY OF WATERLOO**

# The XOR problem



Original $x$ space — Learned $h$ space

- $X = \{ (0,0), (0,1), (1,0), (1,1) \}$, $y = \{-1, 1, 1, -1\}$
- $f^*(\mathbf{x}) = \text{sign}[ f_{xor} - \frac{1}{2} ]$, $f_{xor}(\mathbf{x}) = (x_1 \& -x_2) | (-x_1 \& x_2)$

UNIVERSITY OF WATERLOO

# No separating hyperplane

$$y(\langle \mathbf{w}, \mathbf{x} \rangle + b) > 0$$

- $\mathbf{x}_1 = (0,0)$, $y_1 = -1$  ➔  b < 0
- $\mathbf{x}_2 = (0,1)$, $y_2 = 1$  ➔  $w_2 + b > 0$  ⎤
- $\mathbf{x}_3 = (1,0)$, $y_3 = 1$  ➔  $w_1 + b > 0$  ⎦  $(w_1 + w_2 + b) + b > 0$
- $\mathbf{x}_4 = (1,1)$, $y_4 = -1$  ➔  $w_1 + w_2 + b < 0$

- Contradiction!

Ex. what happens if we run perceptron/winnow on this example?

- [At least one of the blue or green inequalities has to be strict]
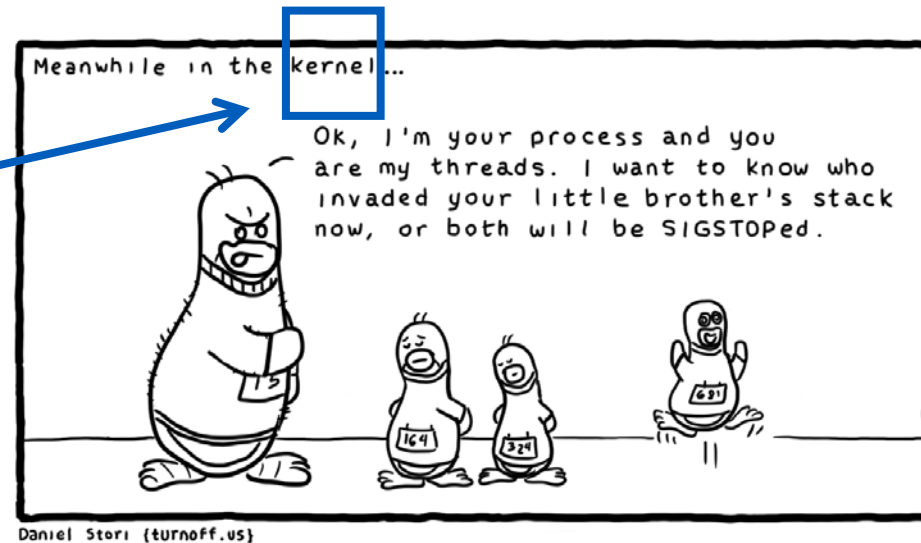
UNIVERSITY OF
WATERLOO

# Fixing the problem

- Our model (hyperplanes) underfits the data (xor func)

- Fix representation, richer model
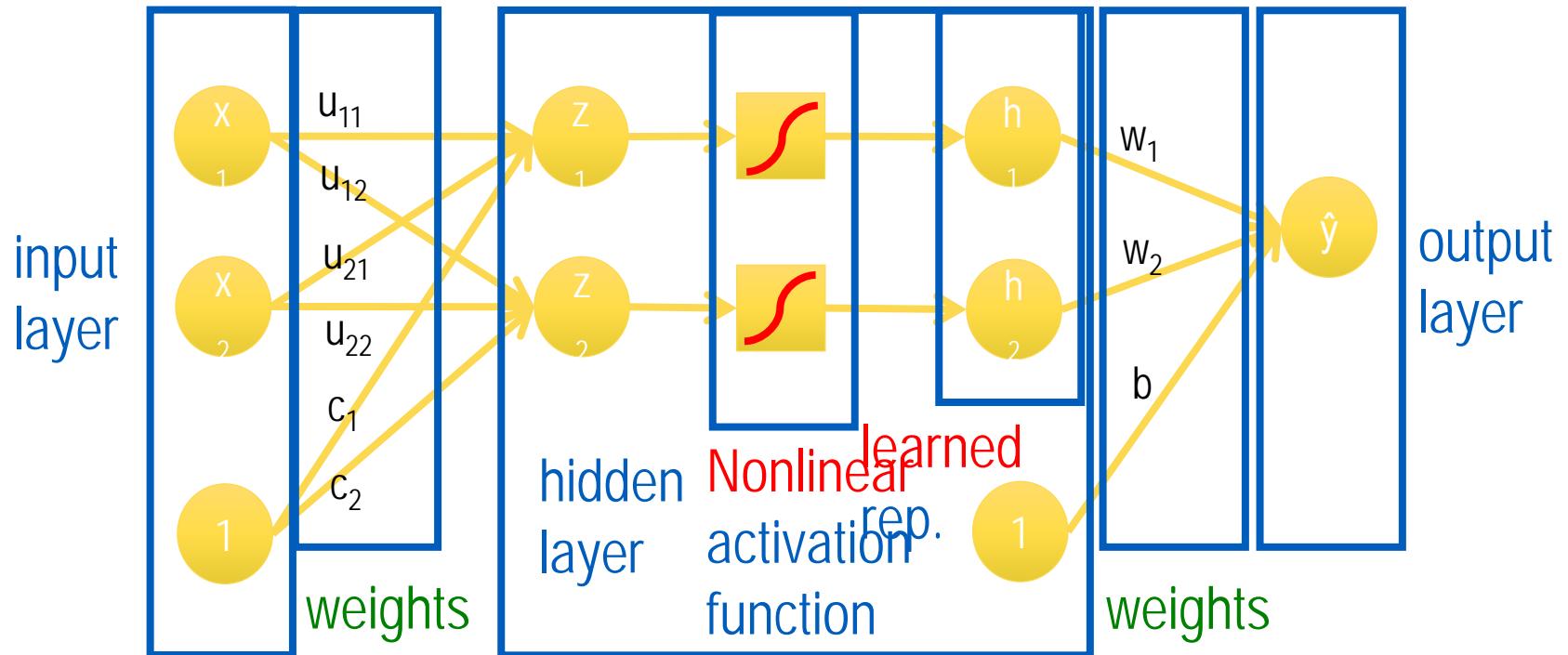
**||**

- Fix model, richer representation



Meanwhile in the kernel...

Ok, I'm your process and you are my threads. I want to know who invaded your little brother's stack now, or both will be SIGSTOPed.

Daniel Stori {turnoff.us}

- NN: still use hyperplane, but learn representation simultaneously

9/19/17                           Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Outline

- Failure of Perceptron

- Neural Network

- Backpropagation

- Universal Approximator

9/19/17                          Yao-Liang Yu

# Two-layer perceptron



input layer

output layer

hidden layer

Nonlinear activation function

learned rep.

weights

weights

$$\mathbf{z} = U\mathbf{x} + \mathbf{c}$$  1st linear layer

makes all the difference!  $$\mathbf{h} = f(\mathbf{z})$$  nonlinear transform
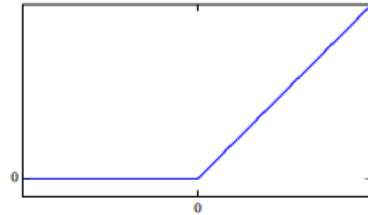
2nd linear layer  $$\hat{y} = \langle \mathbf{h}, \mathbf{w} \rangle + b$$

9/19/17

Yao-Liang Yu

UNIVERSITY OF **WATERLOO**

# Does it work?

$$U = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 2 \\ -4 \end{bmatrix} \quad b = -1$$
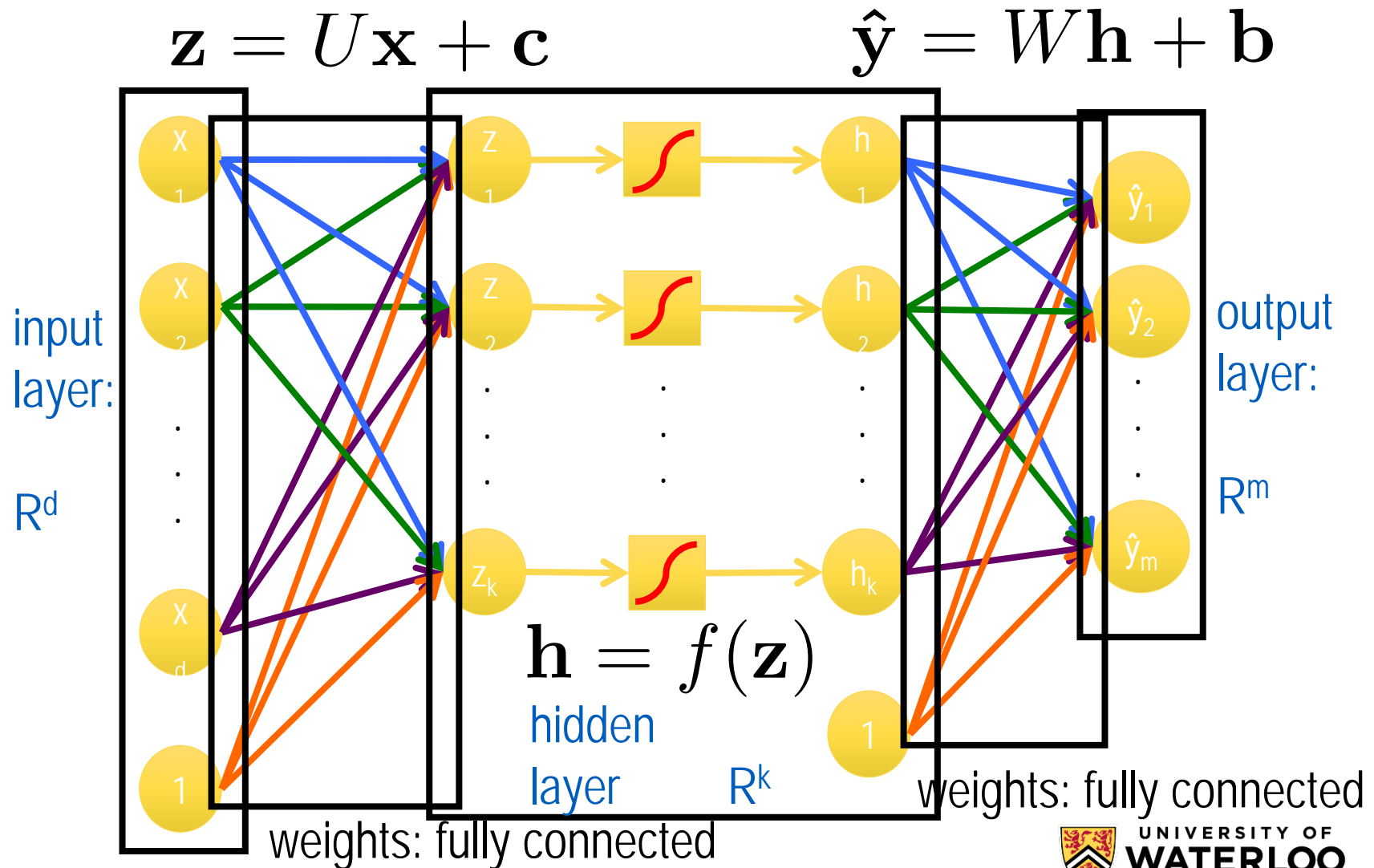
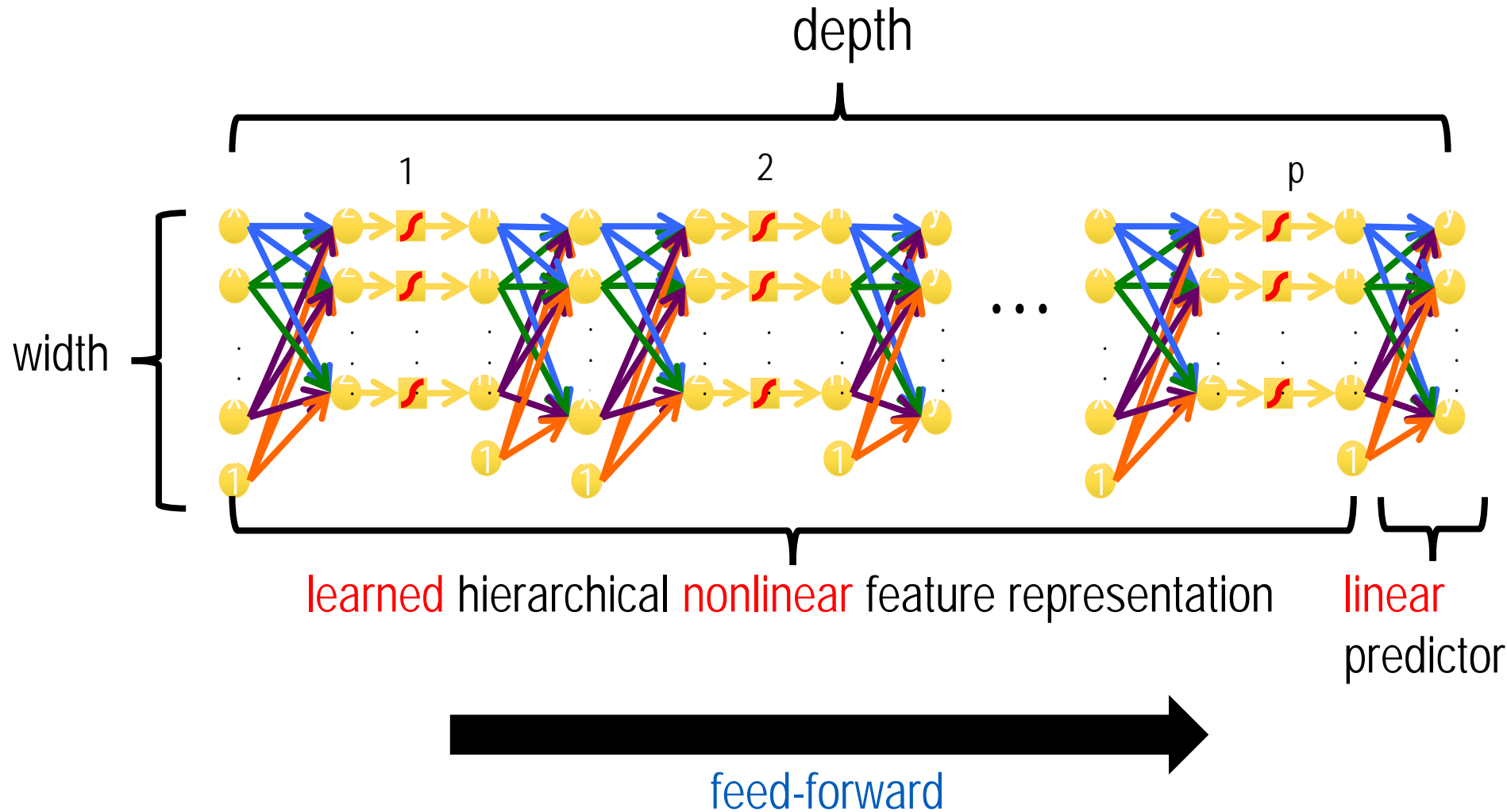Rectified Linear Unit (ReLU) $\qquad f(t) = t_+ := \max(t, 0)$

- $\mathbf{x}_1 = (0,0)$, $y_1 = -1$ → $\mathbf{z}_1 = (0,-1)$, $\mathbf{h}_1 = (0,0)$ → $\hat{y}_1 = -1$ ✔
- $\mathbf{x}_2 = (0,1)$, $y_2 = 1$ → $\mathbf{z}_2 = (1,0)$, $\mathbf{h}_2 = (1,0)$ → $\hat{y}_2 = 1$ ✔
- $\mathbf{x}_3 = (1,0)$, $y_3 = 1$ → $\mathbf{z}_3 = (1,0)$, $\mathbf{h}_3 = (1,0)$ → $\hat{y}_3 = 1$ ✔
- $\mathbf{x}_4 = (1,1)$, $y_4 = -1$ → $\mathbf{z}_4 = (2,1)$, $\mathbf{h}_4 = (2,1)$ → $\hat{y}_4 = -1$ ✔

9/19/17

Yao-Liang Yu

UNIVERSITY OF WATERLOO

# Multi-layer perceptron



$$\mathbf{z} = U\mathbf{x} + \mathbf{c}$$

$$\hat{\mathbf{y}} = W\mathbf{h} + \mathbf{b}$$

input layer:

$R^d$

$$\mathbf{h} = f(\mathbf{z})$$

hidden layer $R^k$

output layer:

$R^m$

weights: fully connected

weights: fully connected

9/19/17

Yao-Liang Yu

UNIVERSITY OF WATERLOO

# Multi-layer perceptron (stacked)



depth

1      2      p

width

learned hierarchical nonlinear feature representation    linear predictor

feed-forward

UNIVERSITY OF
WATERLOO

# Activation function

- Sigmoid

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

- Tanh

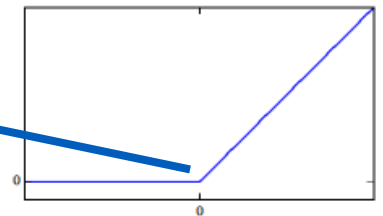$$f(t) = \tanh(t) = \frac{e^t - e^{-t}}{e^t + e^{-t}}$$
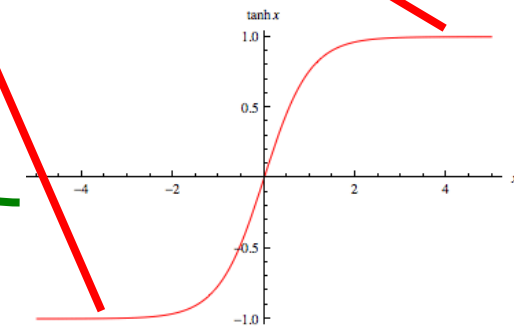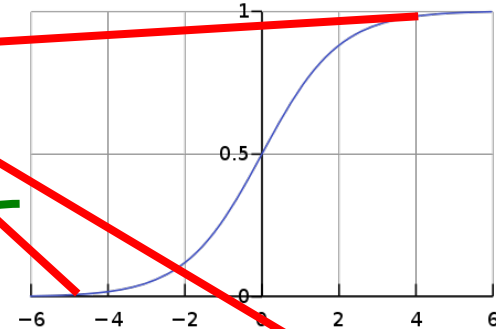
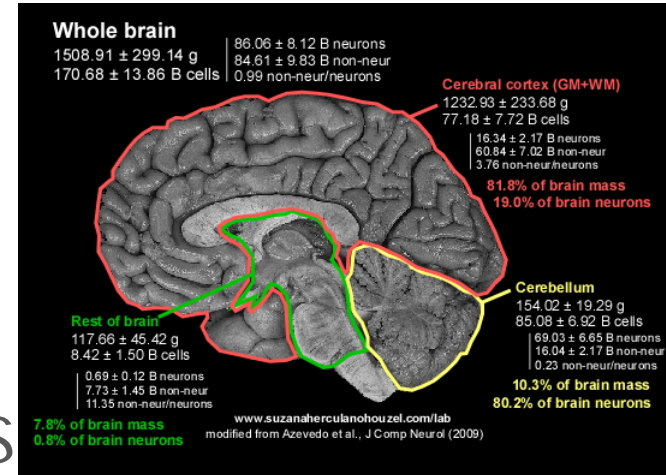- Rectified Linear

$$f(t) = t_+ := \max(t, 0)$$

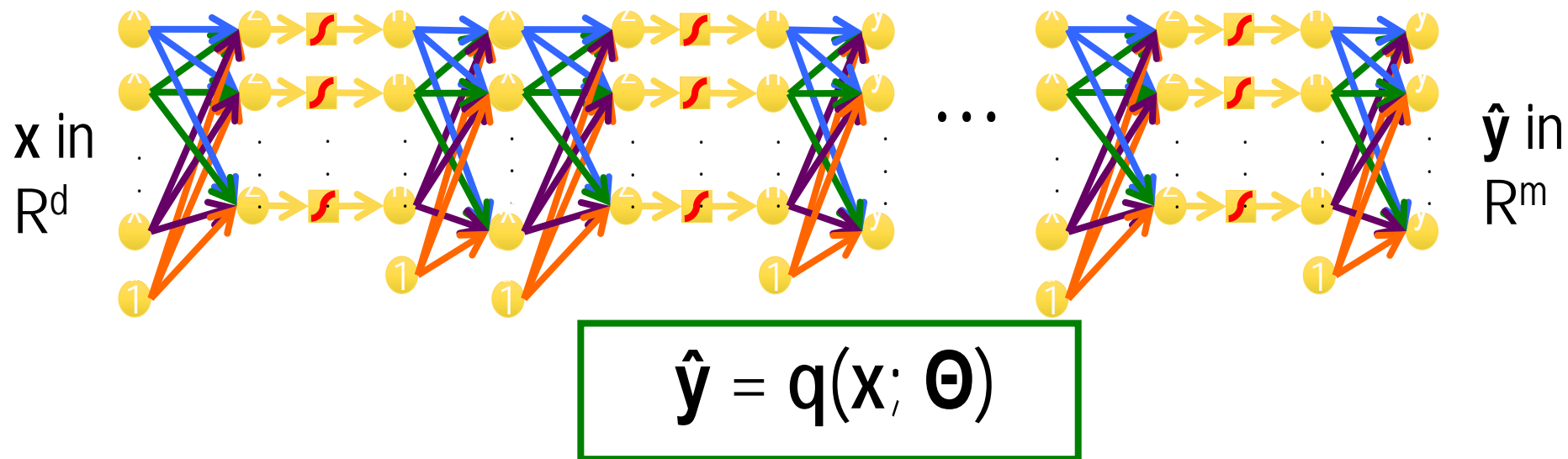饱和 saturate

smooth

nonsmooth

UNIVERSITY OF
**WATERLOO**

# Underfitting vs. Overfitting

- Linear predictor (perceptron / winnow / linear regression) underfits

- NNs learn *hierarchical nonlinear* feature jointly with linear predictor
  - may overfit
  - tons of heuristics (some later)



- Size of NN: # of weights/connections
  - ~ 1 billion; human brain $10^6$ billions (estimated)

# Weights training



**x** in R<sup>d</sup>

**ŷ** in R<sup>m</sup>

$$\hat{\mathbf{y}} = q(\mathbf{x}; \Theta)$$

- Need a loss $\ell$ to measure diff. between pred **ŷ** and truth **y**
  - E.g., $(\hat{\mathbf{y}}-\mathbf{y})^2$; more later

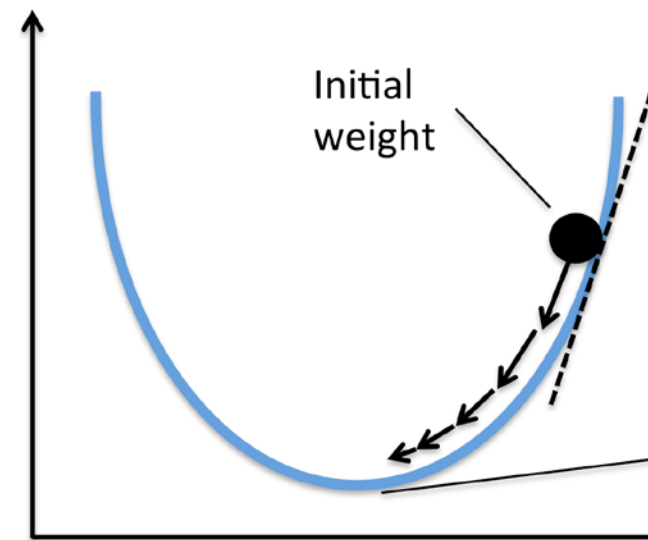- Need a training set $\{(\mathbf{x}_1,\mathbf{y}_1), \ldots, (\mathbf{x}_n,\mathbf{y}_n)\}$ to train weights $\Theta$

UNIVERSITY OF
**WATERLOO**

# Gradient Descent

$$\underset{\Theta}{\min} \; L(\Theta) := \frac{1}{n} \sum_{i=1}^{n} \ell\big[\mathbf{q}(\mathbf{x}_i; \Theta), \mathbf{y}_i\big]$$

(Generalized) gradient

O(n) !

$$\Theta_{t+1} \leftarrow \Theta_t - \eta_t \nabla L(\Theta_t)$$

Initial weight

Step size (learning rate)
- const., if L is smooth
- diminishing, otherwise
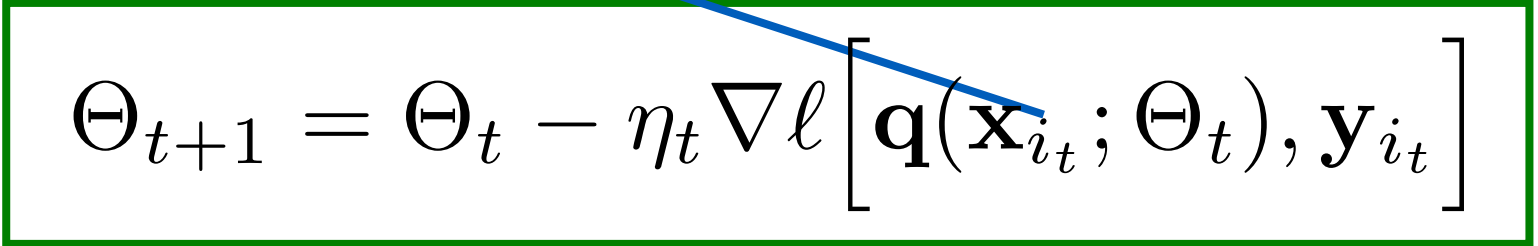
UNIVERSITY OF
**WATERLOO**

# Stochastic Gradient Descent (SGD)

$$\Theta_{t+1} = \Theta_t - \eta_t \cdot \frac{1}{n} \sum_{i=1}^{n} \nabla \ell \Big[ \mathbf{q}(\mathbf{x}_i; \Theta_t), \mathbf{y}_i \Big]$$

average over *n* samples

a random sample suffices

$$\Theta_{t+1} = \Theta_t - \eta_t \nabla \ell \Big[ \mathbf{q}(\mathbf{x}_{i_t}; \Theta_t), \mathbf{y}_{i_t} \Big]$$

- diminishing step size, e.g., 1/sqrt{t} or 1/t
- averaging, momentum, variance-reduction, etc.
- sample w/o replacement; cycle; permute in each pass

# A little history on optimization

- Gradient descent mentioned first in (Cauchy, 1847)

ANALYSE MATHÉMATIQUE. — *Méthode générale pour la résolution des systèmes d'équations simultanées; par* M. AUGUSTIN CAUCHY.

- First rigorous convergence proof (Curry, 1944)

THE METHOD OF STEEPEST DESCENT FOR NON-LINEAR MINIMIZATION PROBLEMS*

BY HASKELL B. CURRY (*Frankford Arsenal*)

- SGD proposed and analyzed (Robbins & Monro, 1951)

A STOCHASTIC APPROXIMATION METHOD[1]

BY HERBERT ROBBINS AND SUTTON MONRO

*University of North Carolina*

Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Herbert Robbins (1915 – 2001)

9/19/17                                        Yao-Liang Yu

# Outline

- Failure of Perceptron

- Neural Network

- Backpropagation

- Universal Approximator

9/19/17     Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**

# Backpropogation

- A fancy name for the chain rule for derivatives:

$$f(x) = g[ h(x) ] \quad \Rightarrow \quad f'(x) = g'[ h(x) ] * h'(x)$$

- Efficiently computes the derivative in NN

- Two passes; complexity = O(size(NN))
  - forward pass: compute function value sequentially
  - backward pass: compute derivative sequentially

9/19/17                                      Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**

# Algorithm

**Require:** Network depth, $l$
**Require:** $W^{(i)}, i \in \{1, \ldots, l\}$, the weight matrices
**Require:** $b^{(i)}, i \in \{1, \ldots, l\}$, the bias parameters ...
**Require:** $x$, the input to process
**Require:** $y$, the target output

$\quad h^{(0)} = x$
$\quad$ **for** $k = 1, \ldots, l$ **do**
$\quad\quad a^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$
$\quad\quad h^{(k)} = f(a^{(k)})$
$\quad$ **end for**
$\quad \hat{y} = h^{(l)}$
$\quad J = L(\hat{y}, y) + \lambda \Omega(\theta)$

Forward pass

f: activation function

$J = L + \lambda\Omega$: training obj.

After the forward computation, compute the gradient on the output layer:
$$g \leftarrow \nabla_{\hat{y}} J = \nabla_{\hat{y}} L(\hat{y}, y)$$
**for** $k = l, l-1, \ldots, 1$ **do**
$\quad$ Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if $f$ is element-wise):
$$g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$$
$\quad$ Compute gradients on weights and biases (including the regularization term, where needed):
$$\nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$$
$$\nabla_{W^{(k)}} J = g \, h^{(k-1)\top} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$$
$\quad$ Propagate the gradients w.r.t. the next lower-level hidden layer's activations:
$$g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)\top} g$$
**end for**

Backward pass

**WATERLOO**

UNIVERSITY OF
WATERLOO

# Outline

- Failure of Perceptron

- Neural Network

- Backpropagation

- Universal Approximator

9/19/17                                    Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**

# Rationals are dense in R

- Any real number can be approximated by some rational number arbitrarily well

- Or in fancy mathematical language

$$\forall r \in \mathbf{R}, \forall \epsilon > 0, \exists s \in \mathbf{Q}, \text{ such that } |r - s| < \epsilon$$

domain of interest     subset in pocket     metric for approx.

UNIVERSITY OF
WATERLOO

# Kolmogorov-Arnold Theorem

Theorem (Kolmogorov, Arnold, Lorentz, Sprecher, ...).
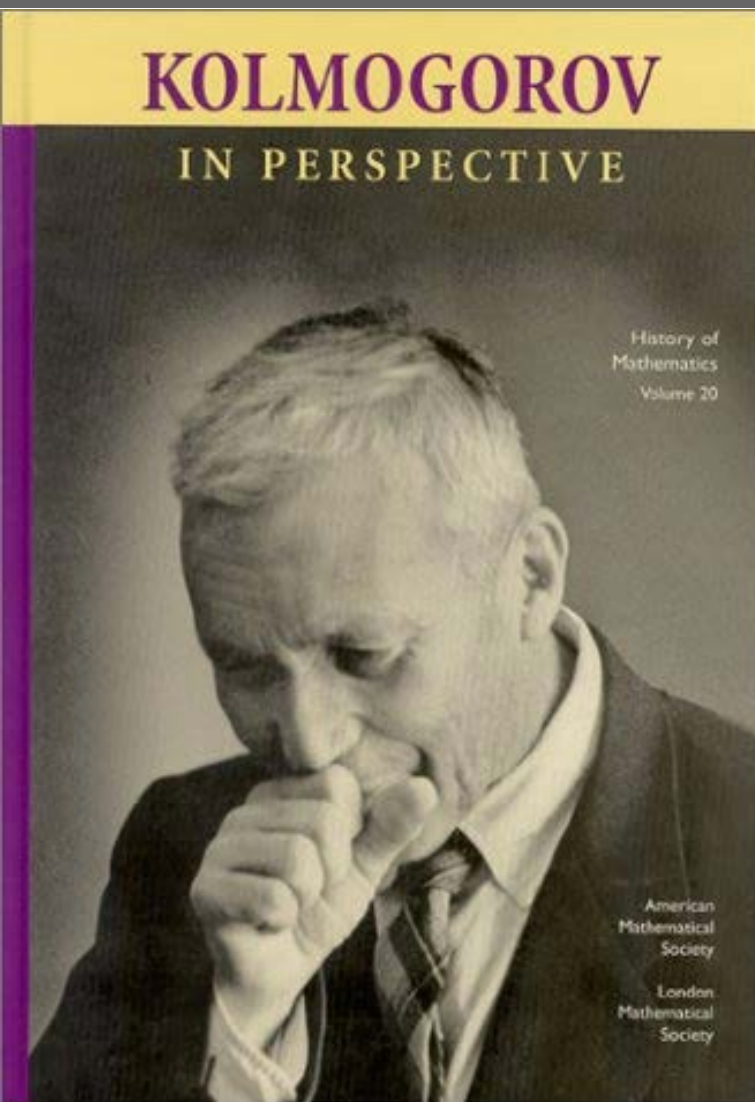Any continuous function $g$: $[0,1]^d \rightarrow R$ can be written
(exactly!) as:
$$g(\mathbf{x}) = \sum_{j=1}^{2d+1} \varphi\Big(\sum_{i=1}^{d} \lambda_i \psi(x_i + \eta j) + j\Big)$$

- Binary addition is the "only" multivariate function!
- Solves Hilbert's 13th problem!
- $\varphi, \psi$ can be constructed from $g$, which is unknown...

Yao-Liang Yu

UNIVERSITY OF
WATERLOO

**KOLMOGOROV**

IN PERSPECTIVE

History of
Mathematics
Volume 20

American
Mathematical
Society

London
Mathematical
Society

ERGEBNISSE DER MATHEMATIK
UND IHRER GRENZGEBIETE

HERAUSGEGEBEN VON DER SCHRIFTLEITUNG
DES
„ZENTRALBLATT FÜR MATHEMATIK"
ZWEITER BAND
——————— 3 ———————

GRUNDBEGRIFFE DER
WAHRSCHEINLICHKEITS-
RECHNUNG

VON

A. KOLMOGOROFF

Foundations of the theory of probability

"Every mathematician believes that he is ahead of the others. The reason none state this belief in public is because they are intelligent people."

BERLIN
VERLAG VON JULIUS SPRINGER
1933

UNIVERSITY OF
**WATERLOO**

# Universal Approximator

**Theorem (Cybenko, Hornik et al., Leshno et al., …).**
Any continuous function $g$: $[0,1]^d \rightarrow$ R can be <span style="color:red">uniformly approximated</span> in arbitrary precision by a <span style="color:red">two-layer</span> NN with an activation function $f$ that
1. is locally bounded
2. has "negligible" closure of discontinuous points
3. is not a polynomial

- conditions are necessary in some sense
- includes (almost) all activation functions in practice

Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Caveat and Remedy

- NNs were praised for being "universal"
  - but shall see later that many kernels are universal as well
  - desirable but perhaps not THE explanation

- May need exponentially many hidden units…

- Increase depth may reduce network size, exponentially!
  - can be a course project, ask for references

UNIVERSITY OF
**WATERLOO**

# Questions?

9/19/17

Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**