Lecture Notes for
**Machine Learning in Python**

Professor Eric Larson
**Town Hall + MLP History**

1

# Class Logistics and Agenda

- Logistics:
  - Next time: Flipped Module on back propagation
- Multi Week Agenda:
  - Today: Neural Networks History, up to 1980
  - Today: Multi-layer Architectures
  - Town Hall, Lab 3 (if time)
  - Flipped: Programming Multi-layer training
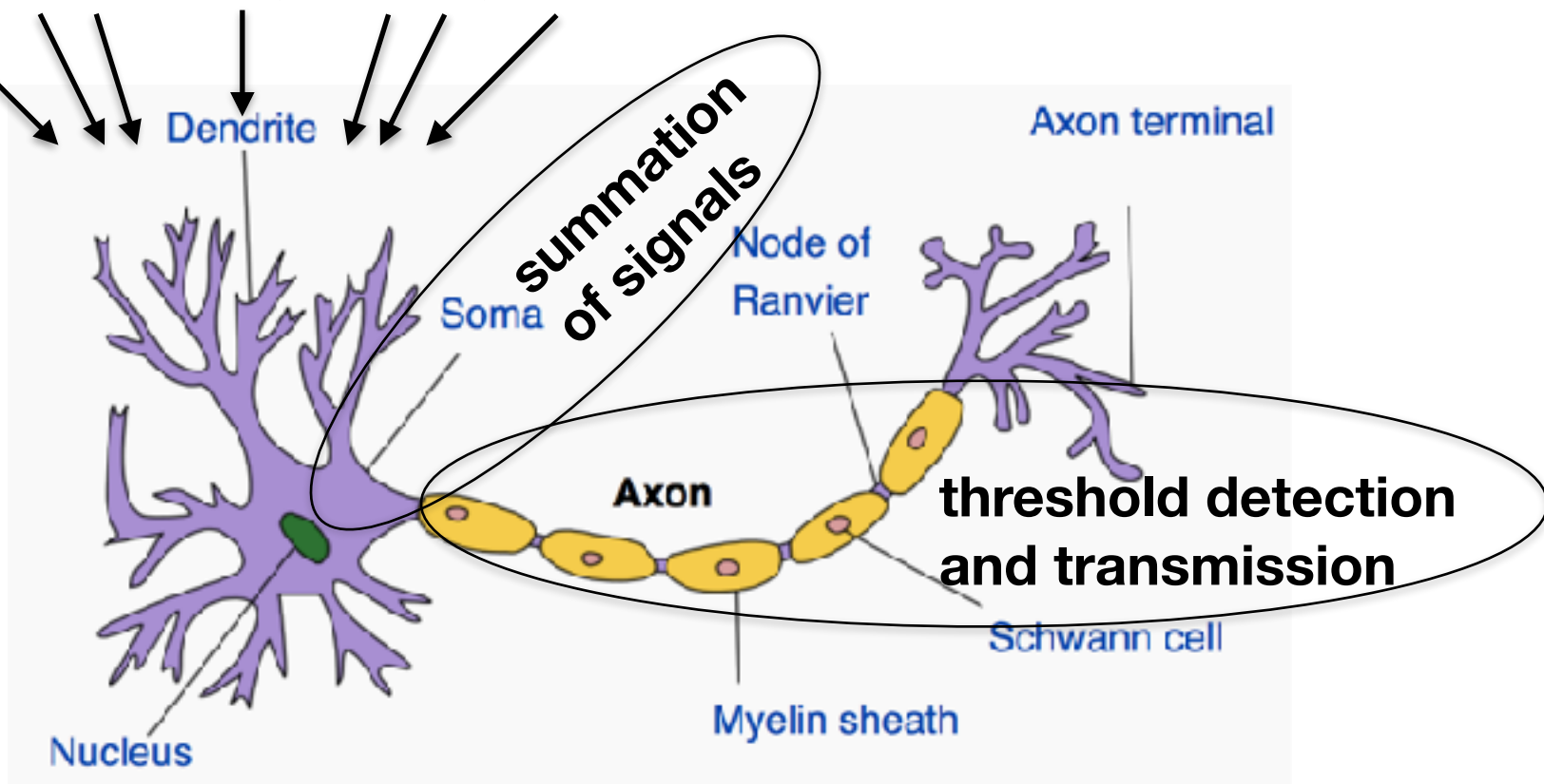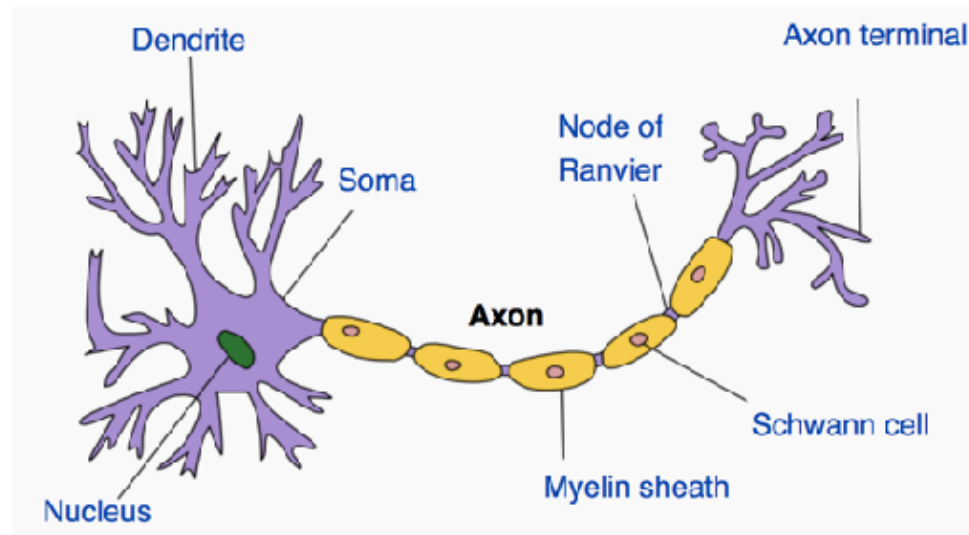
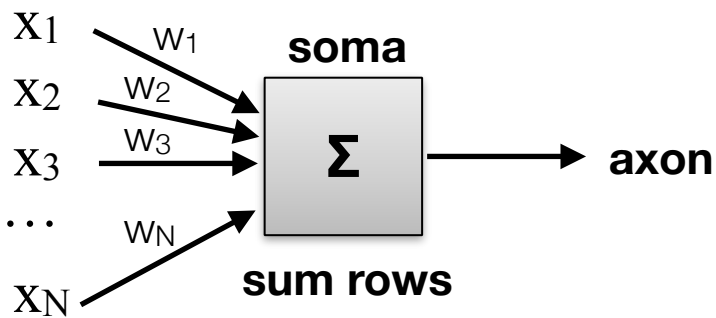# A History of Neural Networks

# Neurons

- From biology to modeling:

**input from neighboring neurons**



summation of signals

**threshold detection and transmission**

**dendrite**

$x_1$ $w_1$

$x_2$ $w_2$

$x_3$ $w_3$

$\cdots$ $w_N$

$x_N$

**soma**

$\Sigma$

**sum rows**

**axon**

**input**
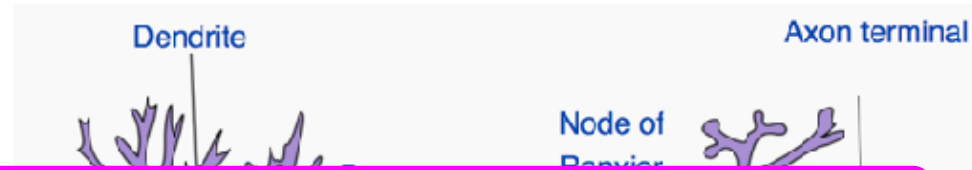
## logic gates of the mind



Warren McCulloch



Walter Pitts

# Neurons

- McCulloch and Pitts, 1943
- Donald Hebb, 1949
  - Hebb's Law: close neurons fire together
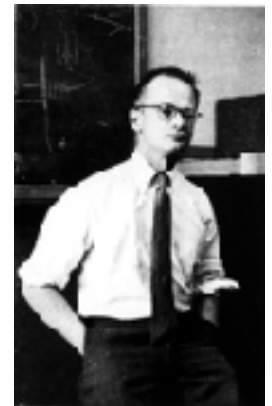  - neurons "learn
  - easier synaptic
  - basis of neura

I as infatuated with the idea of **brainwashing** and controlling minds of others! I also invented a number of **torture procedures** like sensory deprivation and **isolation tanks**—and carried out a number of secret studies on real people!!
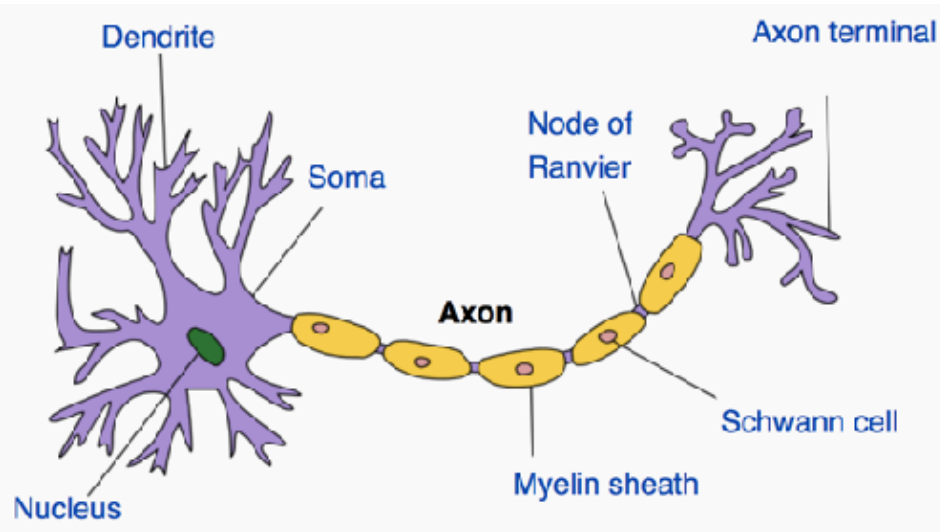
Donald O. Hebb
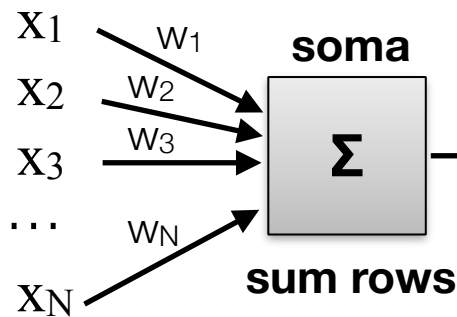
Warren McCulloch

Walter Pitts

# Rosenblatt's perceptron, 1957



Frank Rosenblatt

Dendrite
Soma
Node of Ranvier
Axon terminal
Axon
Schwann cell
Myelin sheath
Nucleus

hard limit

$a = -1 \quad z < 0$
$a = 1 \quad z >= 0$

linear

$a = z$

sigmoid

$a = \dfrac{1}{1 + \exp(-z)}$

**dendrite**

$x_1$ $\quad$ $w_1$

$x_2$ $\quad$ $w_2$

$x_3$ $\quad$ $w_3$

$\ldots$ $\quad$ $w_N$

$x_N$

**soma**

$\Sigma$

**sum rows**

**axon**

$\varphi$

**activation function**

**input**

PERCEPTRON

Perceptron Learning Rule:
~Stochastic Gradient Descent

5 4 3 2 1 0

**dendrite**



**soma**

$\Sigma$

**sum rows**

**axon**

$\varphi$

**activation function**

**input**

$\mathbf{a} = \mathbf{x} +$ concat bias term

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_1 \\ \vdots \\ x_j \\ \vdots \\ x_N \end{bmatrix}^{(i)} \qquad \mathbf{a} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_j \\ \vdots \\ x_N \end{bmatrix}$$

$$\mathbf{z} = \mathbf{W} \cdot \mathbf{a} = \mathbf{W}_{1:N} \cdot \mathbf{x}^{(i)} + \mathbf{b}$$

$$\mathbf{W} = \begin{bmatrix} w_{1,0} \; w_{1,1} \; w_{1,2} \; ... \; w_{1,N} \\ ... \\ w_{S,0} \; w_{S,1} \; w_{S,2} \; ... \; w_{S,N} \end{bmatrix}$$
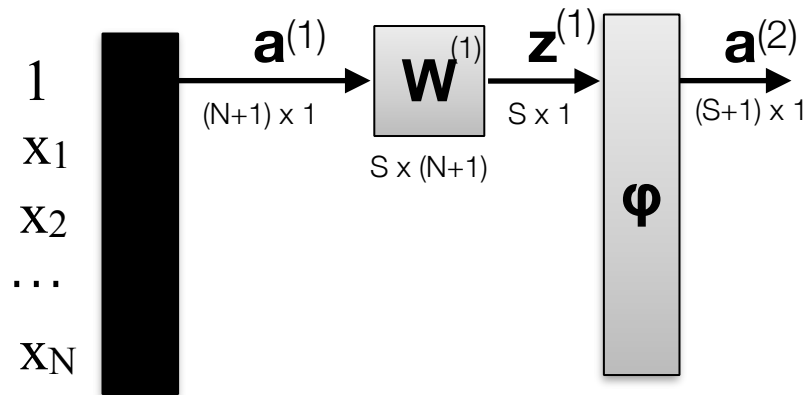
$\mathbf{a}^{(1)}$ $\quad$ $\mathbf{W}^{(1)}$ $\quad$ $\mathbf{z}^{(1)}$ $\quad$ $\varphi$ $\quad$ $\mathbf{a}^{(2)}$

(N+1) x 1 $\qquad$ S x (N+1) $\qquad$ S x 1 $\qquad$ (S+1) x 1

$$[\mathbf{z}^{(1)}]^{(i)} = \mathbf{W}^{(1)} \cdot [\mathbf{a}^{(1)}]^{(i)} = \mathbf{W}^{(1)}_{1:N} \cdot \mathbf{x}^{(i)} + \mathbf{b}^{(1)}$$

$\mathbf{a}^{(next)} = \varphi(\mathbf{z}^{(current)}) +$ concat bias term
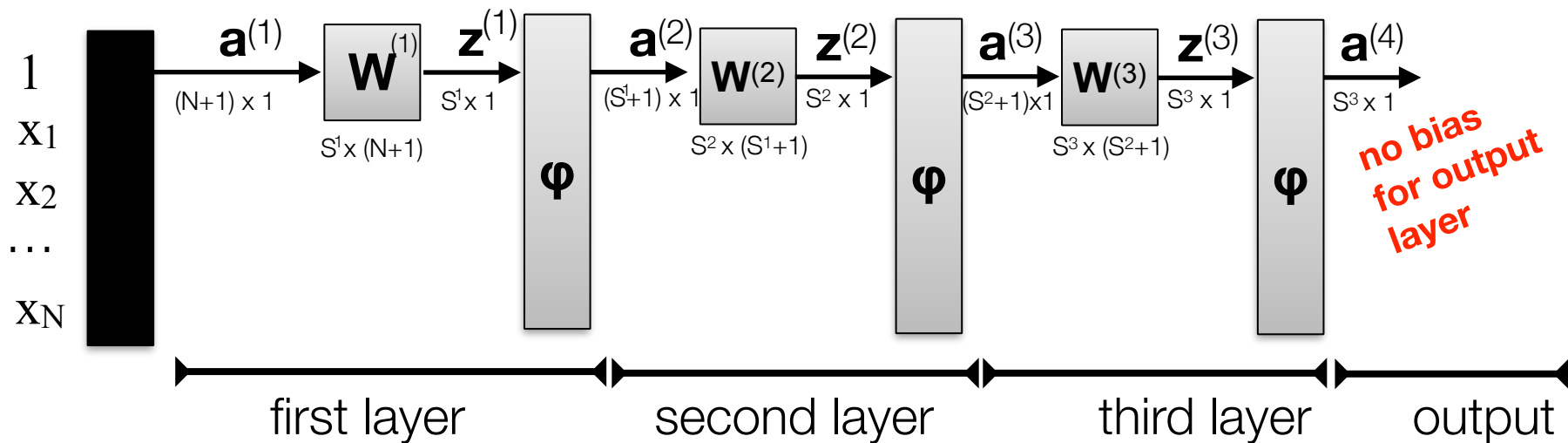
$$\mathbf{a}^{(next)} = \begin{bmatrix} 1 \\ \phi(z_1^{curr}) \\ \vdots \\ \phi(z_N^{curr}) \end{bmatrix} \rightarrow \mathbf{a}^{(L)} = \begin{bmatrix} 1 \\ \phi(z_1^{L-1}) \\ \vdots \\ \phi(z_N^{L-1}) \end{bmatrix}$$

$\mathbf{x}^{(i)}$ One row from Table data becomes input column to model

notation adapted from *Neural Network Design*, Hagan, Demuth, Beale, and De Jesus

**9**

$$\mathbf{a}^{(L+1)} = \boldsymbol{\phi}(\mathbf{z}^{(L)}) + \text{concat bias term}$$

$$\mathbf{a}^{(4)} \quad \text{rows=unique classes}$$

$$\mathbf{z}^{(L)} = \mathbf{W}^{(L)}\mathbf{a}^{(L)}$$

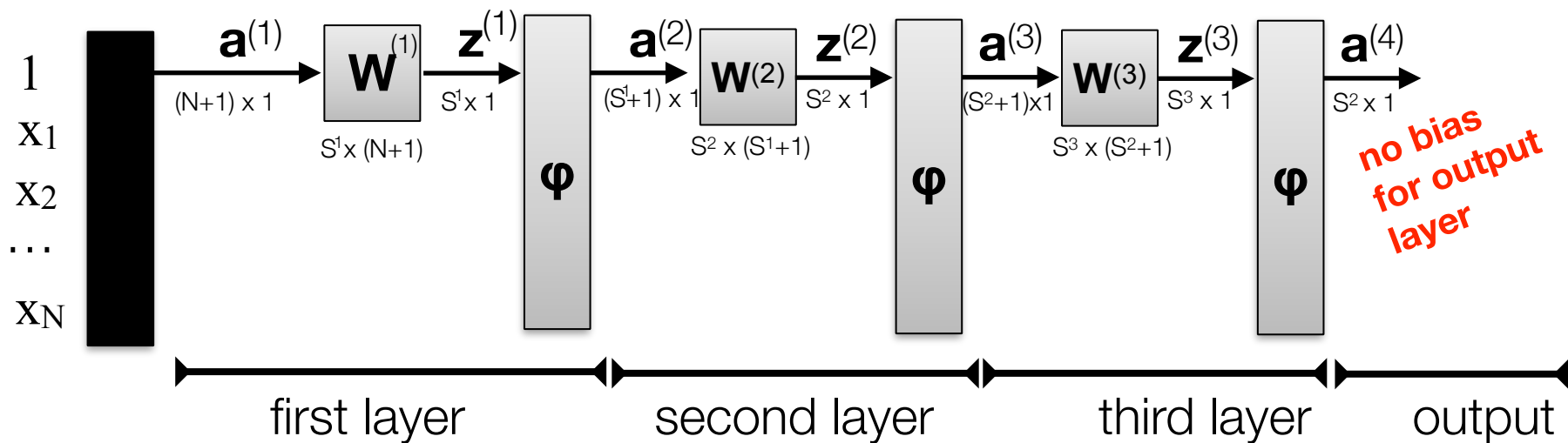$$\mathbf{z}^{(L)} = \mathbf{W}^{(L)}\boldsymbol{\phi}(\mathbf{z}^{(L-1)})$$

$$\mathbf{W}^{(L)} = \begin{bmatrix} w^{(L)}_{1,0} & w^{(L)}_{1,1} & w^{(L)}_{1,2} & \ldots & w^{(L)}_{1,S^{L-1}} \\ & & \ldots & & \\ w^{(L)}_{S^L,0} & w^{(L)}_{S^L,1} & \ldots & & w^{(L)}_{S^L,S^{L-1}} \end{bmatrix}$$

$S^L \times (S^{L-1}+1)$

$$\mathbf{z}^{(L)} = \mathbf{W}^{(L)}\boldsymbol{\phi}(\mathbf{W}^{(L-1)}\mathbf{a}^{(L-1)})$$

**10**

# Multiple layers notation



- **Self test**: How many parameters need to be trained in the above network?
  - A. $[(N+1) \times S^1] \;\; + \;\; [(S^1+1) \times S^2] \;\; + \;\; [(S^2+1) \times S^3]$
  - B. $|\mathbf{W}^{(1)}| + |\mathbf{W}^{(2)}| + |\mathbf{W}^{(3)}|$
  - C. can't determine from diagram
  - D. it depends on the sizes of intermediate variables, $\mathbf{z}^{(i)}$

notation adapted from *Neural Network Design*, Hagan, Demuth, Beale, and De Jesus  **11**

$$z^{(L)}=W^{(L)}a^{(L)}$$

$$[z^{(L)}]^{(i)}=W^{(L)}[a^{(L)}]^{(i)}$$

$$Z^{(L)}=W^{(L)}A^{(L)}$$

$$Z^{(L)}=W^{(L)}\phi(W^{(L-1)}A^{(L-1)})$$

# Training Neural Network Architectures



MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED
THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.

$$\sum_i^M (\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)})^2$$

Need objective Function, minimize MSE

$$J(\mathbf{W}) = ||\mathbf{Y} - \hat{\mathbf{Y}}||^2$$

where ground truth $\mathbf{y}^{(i)}$ is one-hot encoded!

$$\begin{bmatrix} y_1 \\ ... \\ y_C \end{bmatrix}^{(i)}$$

$$\begin{bmatrix} \begin{bmatrix} y_1 \\ ... \\ y_C \end{bmatrix}^{(1)} \quad ... \quad \begin{bmatrix} y_1 \\ ... \\ y_C \end{bmatrix}^{(M)} \end{bmatrix} = \mathbf{Y}$$

# Simple Architectures

- Rosenblatt's perceptron, 1957



**Self Test** - If this is a binary classification problem, how large is *S*, the length of yhat?

      A. Can't determine

      B. 2

      C. 1

      D. N

# Simple Architectures

- Adaline network, Widrow and Hoff, 1960



linear

phi

Marcian "Ted" Hoff

Bernard Widrow

Objective Function, minimize MSE $J(\mathbf{W})=||\mathbf{Y}-\overset{\triangle}{\mathbf{Y}}||^2$

New objective function becomes: $J(\mathbf{W})=||\mathbf{Y}-\mathbf{W}\cdot\mathbf{A}||^2$

Need gradient $\nabla J(\mathbf{W})$ for update equation $\mathbf{W} \leftarrow \mathbf{W} + \eta \nabla J(\mathbf{W})$

We have been using the **Widrow-Hoff Learning Rule**

# Simple Architectures

- Adaline network, Widrow and Hoff, 1960



Marcian "Ted" Hoff

Bernard Widrow

linear

`phi`

**a** (N+1) x 1 → **W** S x (N+1) → **z** S x1 → **φ** → **yhat** S x 1

need gradient $\nabla J(\mathbf{W})$ for update equation $\mathbf{W} \leftarrow \mathbf{W} + \eta \nabla J(\mathbf{W})$

For case S=1, $\mathbf{W}$ has only one row, $\mathbf{w}$
this is just **linear regression…**

$$\mathbf{w} \leftarrow \mathbf{w} + \eta[\mathbf{X} * (\mathbf{y} - \hat{\mathbf{y}})]$$

# Simple Architectures

- Modern Perceptron network

**a** $(N+1) \times 1$ → **W** $S \times (N+1)$ → **z** $S \times 1$ → **φ** → **yhat** $S \times 1$

sigmoid

$$phi = \frac{1}{1+exp(-z)}$$

need gradient $\nabla J(\mathbf{W})$ for update equation $\mathbf{W} \leftarrow \mathbf{W} + \eta \nabla J(\mathbf{W})$

For case S=1, this is just **logistic regression…** and **we have already solved this**!

$$\mathbf{w} \leftarrow \mathbf{w} + \eta[\mathbf{X} * (\mathbf{y} - \mathbf{g(x)})]$$

## What happens when S > 1 ?

$$\mathbf{yhat}^{(i)} = \begin{bmatrix} \varphi(_{row=1}\mathbf{w}\cdot\mathbf{x}^{(i)}) \\ \varphi(_{row=2}\mathbf{w}\cdot\mathbf{x}^{(i)}) \\ \dots \\ \varphi(_{row=S}\mathbf{w}\cdot\mathbf{x}^{(i)}) \end{bmatrix}$$

one hot

which is one-versus-all!

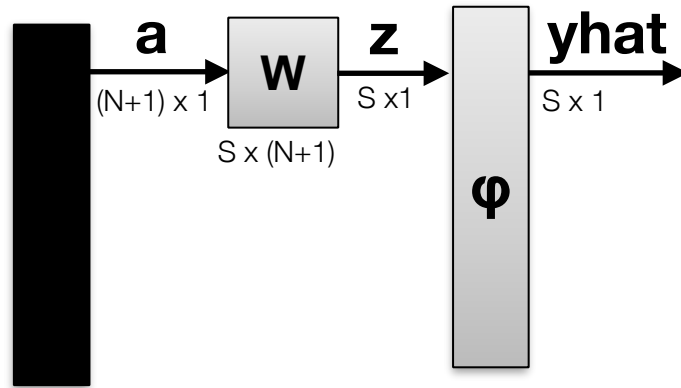$$\begin{bmatrix} \begin{bmatrix} \phi(z_1) \\ \dots \\ \phi(z_S) \end{bmatrix}^{(1)} & \dots & \begin{bmatrix} \phi(z_1) \\ \dots \\ \phi(z_S) \end{bmatrix}^{(M)} \end{bmatrix} = \overset{\triangle}{\mathbf{Y}}$$

$$\begin{bmatrix} \begin{bmatrix} y_1 \\ \dots \\ y_S \end{bmatrix}^{(1)} & \dots & \begin{bmatrix} y_1 \\ \dots \\ y_S \end{bmatrix}^{(M)} \end{bmatrix} = \mathbf{Y}$$

$$J(_1\mathbf{w}) = \sum_{i=1} [y_1^{(i)} - \varphi(_1\mathbf{w}\cdot\mathbf{x}^{(i)})]^2$$

$$J(_2\mathbf{w}) = \sum_{i=1} [y_2^{(i)} - \varphi(_2\mathbf{w}\cdot\mathbf{x}^{(i)})]^2$$

$$\dots$$

$$J(\mathbf{W}) = ||\mathbf{Y} - \overset{\triangle}{\mathbf{Y}}||^2$$

$$J(_S\mathbf{w}) = \sum_{i=1} [y_S^{(i)} - \varphi(_S\mathbf{w}\cdot\mathbf{x}^{(i)})]^2$$

Each target class in $\mathbf{Y}$ can be independently optimized

# Simple Architectures: Summary

- Adaline network, Widrow and Hoff, 1960
  - linear regression
- Perceptron
  - *with sigmoid*: logistic regression
- One-versus-all implementation is the same as having $\mathbf{w}_{class}$ be rows of weight matrix, **W**
  - works in adaline
  - works in logistic regression

these networks were created in the 50's and 60's but were abandoned
**why were they not used?**

# The Rosenblatt-Widrow-Hoff Dilemma

- 1960's: Rosenblatt got into a public academic argument with Marvin Minsky and Seymour Papert

  "Given an elementary α-perceptron, a stimulus world W, and any classification C(W) for which a solution exists; let all stimuli in W occur in any sequence, provided that each stimulus must reoccur in finite time; then beginning from an arbitrary initial state, an error correction procedure will always yield a solution to C(W) in finite time…"

- Minsky and Papert publish limitations paper, 1969:



TED  Ideas worth spreading          WATCH      DISCOVER      ATTEND      PARTICIPATE

Marvin Minsky:

**Health and the human mind**

TED2003 · 13:33 · Filmed Feb 2003

21 subtitle languages

View interactive transcript

# More Advanced Architectures: history

- 1986: *Rumelhart*, *Hinton*, and *Williams* popularize gradient calculation for multi-layer network
    - *technically* introduced by Werbos in 1982
- **difference**: Rumelhart *et al.* validated ideas with a computer
- until this point no one could train a multiple layer network consistently
- algorithm is popularly called **Back-Propagation**
- wins pattern recognition prize in 1993, becomes de-facto machine learning algorithm until: SVMs and Random Forests in ~2004
- would eventually see a resurgence for its ability to train algorithms for Deep Learning applications: **Hinton is widely considered the founder of deep learning**
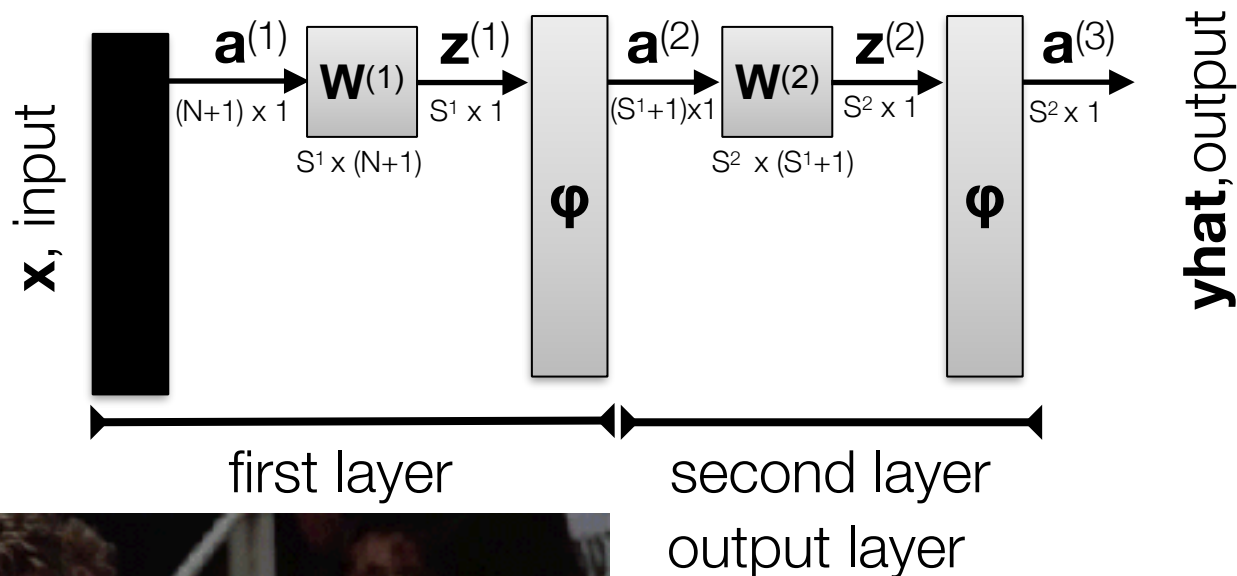
David Rumelhart

1942-2011

Geoffrey Hinton

# More Advanced Architectures: MLP

- The multi-layer perceptron (MLP):
  - two layers shown, but could be arbitrarily many layers
  - algorithm is agnostic to number of layers (*kinda*)

each row of **yhat** is no longer independent of the rows in **W** so we cannot optimize using one versus all!!!

$\mathbf{x}$, input

$\mathbf{a}^{(1)}$ $(N+1) \times 1$

$\mathbf{W}^{(1)}$ $S^1 \times (N+1)$

$\mathbf{z}^{(1)}$ $S^1 \times 1$

$\boldsymbol{\varphi}$

$\mathbf{a}^{(2)}$ $(S^1+1) \times 1$

$\mathbf{W}^{(2)}$ $S^2 \times (S^1+1)$

$\mathbf{z}^{(2)}$ $S^2 \times 1$

$\boldsymbol{\varphi}$

$\mathbf{a}^{(3)}$ $S^2 \times 1$

**yhat**, output

first layer

second layer

output layer

$$\mathbf{yhat}^{(i)} = \begin{bmatrix} \varphi\left(_{\mathrm{row}=1}\mathbf{w}^{(2)} \cdot \varphi(\mathbf{W}^{(1)}\mathbf{a}^{(1)})\right) \\ \cdots \\ \varphi\left(_{\mathrm{row}=S}\mathbf{w}^{(2)} \cdot \varphi(\mathbf{W}^{(1)}\mathbf{a}^{(1)})\right) \end{bmatrix}$$

one hot

SWEEP THE LEG.

23

- Steps:
  - propagate weights forward
  - calculate gradient at final layer
  - back propagate gradient for each layer
    - via recurrence relation
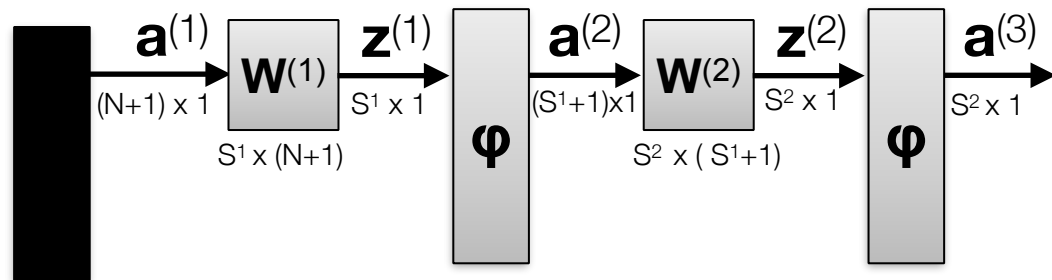


$$J(\mathbf{W}) = \|\mathbf{Y} - \overset{\triangle}{\mathbf{Y}}\|^2$$

$$w_{i,j}^{(l)} \leftarrow w_{i,j}^{(l)} - \eta \frac{\partial J(\mathbf{W})}{\partial w_{i,j}^{(l)}}$$

**x**, input

| **a**$^{(1)}$ | **W**$^{(1)}$ | **z**$^{(1)}$ | **φ** | **a**$^{(2)}$ | **W**$^{(2)}$ | **z**$^{(2)}$ | **φ** | **a**$^{(3)}$ |

$(N+1) \times 1$    $S^1 \times 1$    $(S^1+1) \times 1$    $S^2 \times 1$    $S^2 \times 1$

$S^1 \times (N+1)$    $S^2 \times (S^1+1)$

**yhat**, output

first layer

second layer
output layer

24

# Back propagation

$$J(\mathbf{W}) = \sum_{k}^{M} (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$



$$w_{i,j}^{(l)} \leftarrow w_{i,j}^{(l)} - \eta \cdot \frac{\partial J(\mathbf{W})}{\partial w_{i,j}^{(l)}}$$

use chain rule:

$$\frac{\partial J(\mathbf{W})}{\partial w_{i,j}^{(l)}} = \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} \frac{\partial \mathbf{z}^{(l)}}{\partial w_{i,j}^{(l)}}$$

Solve this in explainer video for next in class assignment!

# Lab 3, Town Hall (if time)

- thanks! **Next time is Flipped Assignment!!!**

**More help on neural networks to prepare for next time:**

Sebastian Raschka

https://github.com/rasbt/python-machine-learning-book/blob/master/code/ch12/ch12.ipynb

Martin Hagan

https://www.google.com/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=0ahUKEwioprvn27fPAhWMx4MKHYbwDlwQFggeMAA&url=http%3A%2F%2Fhagan.okstate.edu%2FNNDesign.pdf&usg=AFQjCNG5YbM4xSMm6K5HNsG-4Q8TvOu_Lw&sig2=bgT3k-5ZDDTPZ07Qu8Oreg

Michael Nielsen

http://neuralnetworksanddeeplearning.com