

# Lecture Notes for **Machine Learning in Python**

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

# 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

# Class Overview, by topic

Table Data  
Visualization

Numpy, Pandas, Seaborn  
Overviews with some in-depth discussion

Dimension  
Reduction and  
Image Processing

Scikit-learn, Scikit Image,  
Intuition only, Some mathematics

Linear and  
Logistic  
Regression

Numpy, Recreate API for Scikit-learn  
Detailed mathematics for simple optimization  
intuition for advanced optimization

Neural Networks  
and Back Prop.

Numpy  
Detailed mathematics for NN operations

Wide and Deep  
Networks

Convolutional  
Networks

Recurrent  
Networks

Keras, Tensorflow  
Intuition, Detailed implement.

Ethics in  
Language Models

ConceptNet  
Case studies

# A History of Neural Networks

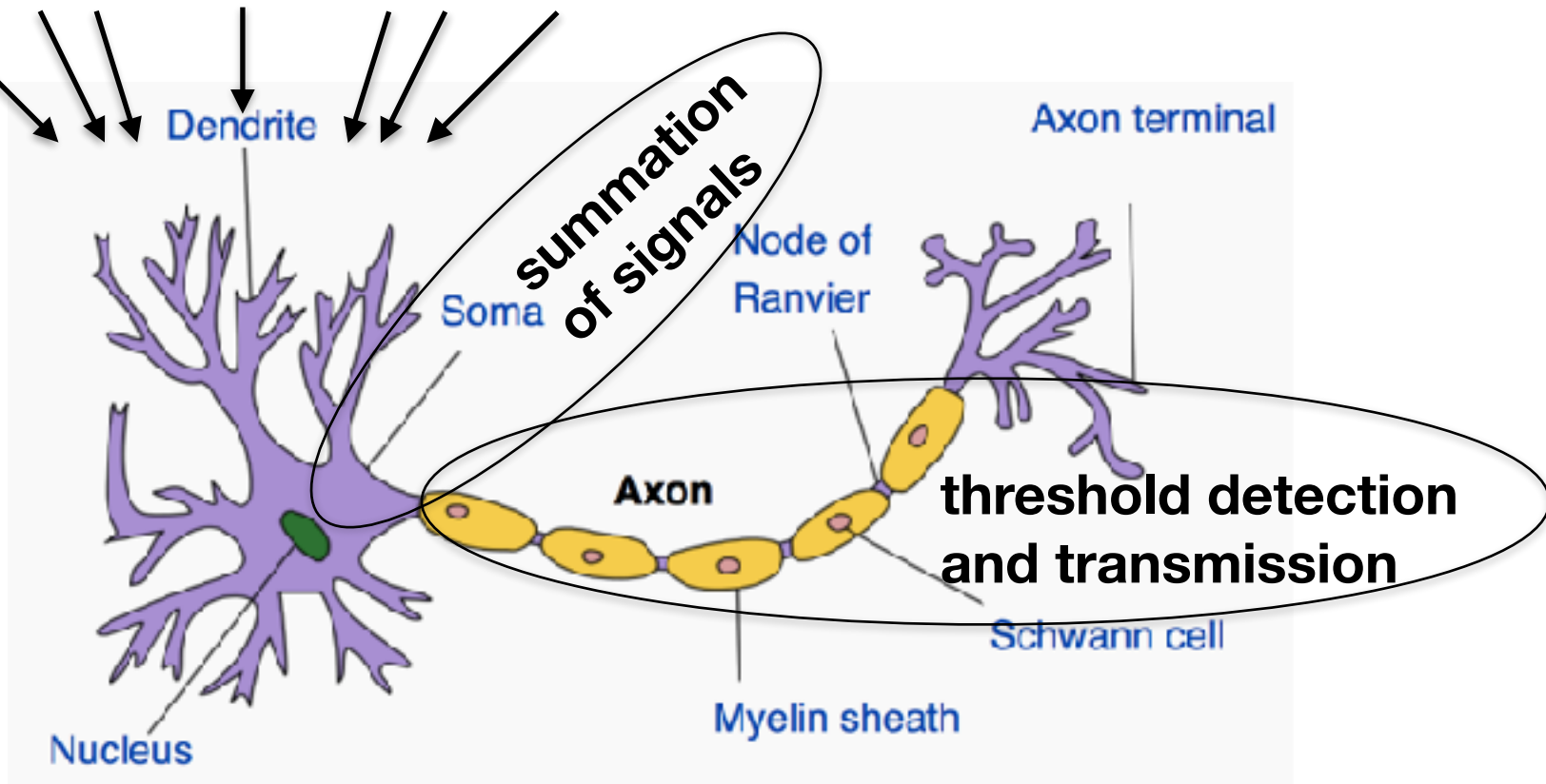


Machine Learning 101

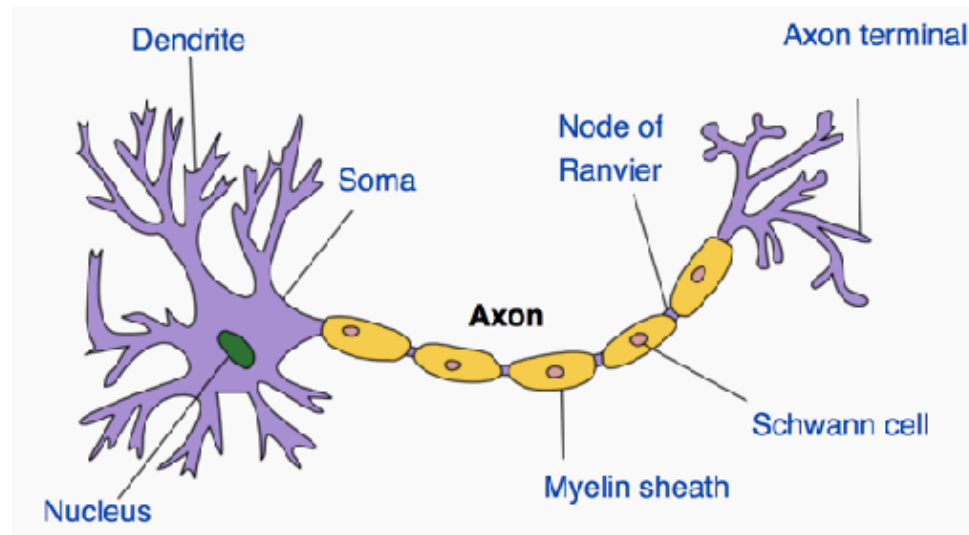
# Neurons

- From biology to modeling:

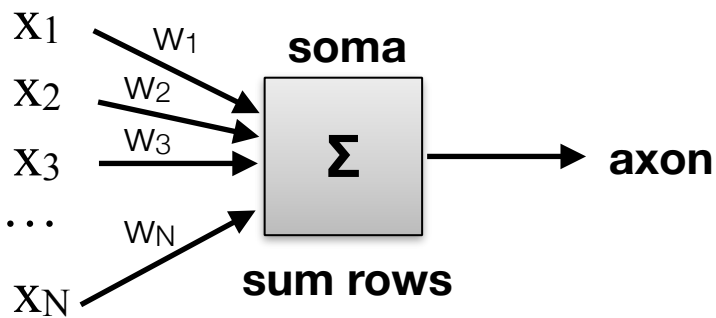
**input from neighboring neurons**



# McCulloch and Pitts, 1943



**dendrite**



**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 neural



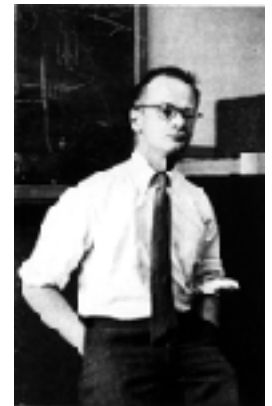
I was infatuated with the idea of **brainwashing** and controlling minds of others! I also invented a number of **torture procedures** like sensory deprivation and **isolation tanks**—and carried out a number of secret studies on real people!!



Donald O. Hebb



Warren McCulloch

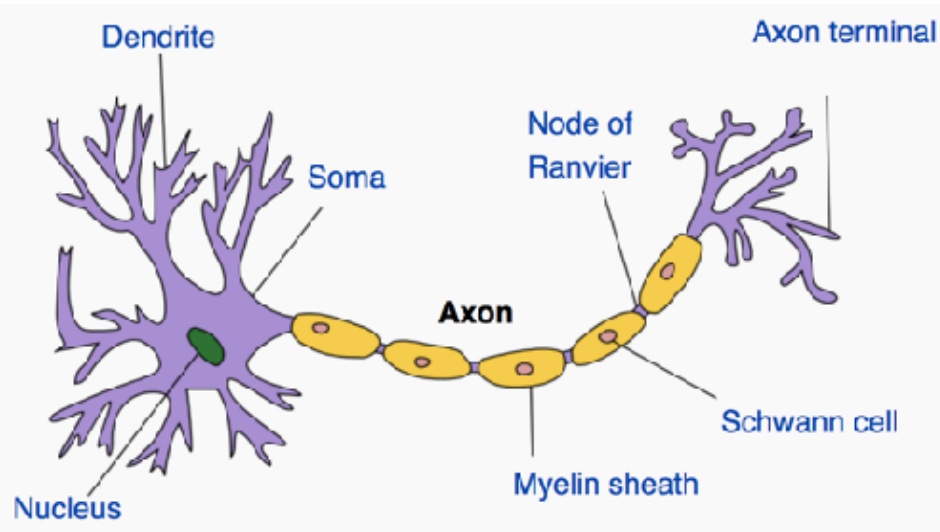


Walter Pitts

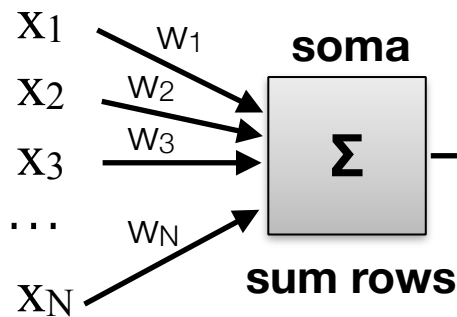
# Rosenblatt's perceptron, 1957



Frank Rosenblatt

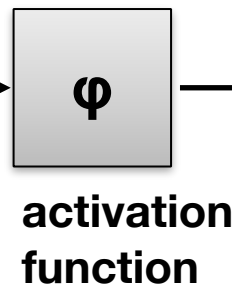


dendrite



input

axon



hard limit



$$\begin{aligned} a &= -1 & z < 0 \\ a &= 1 & z \geq 0 \end{aligned}$$

linear



$$a = z$$

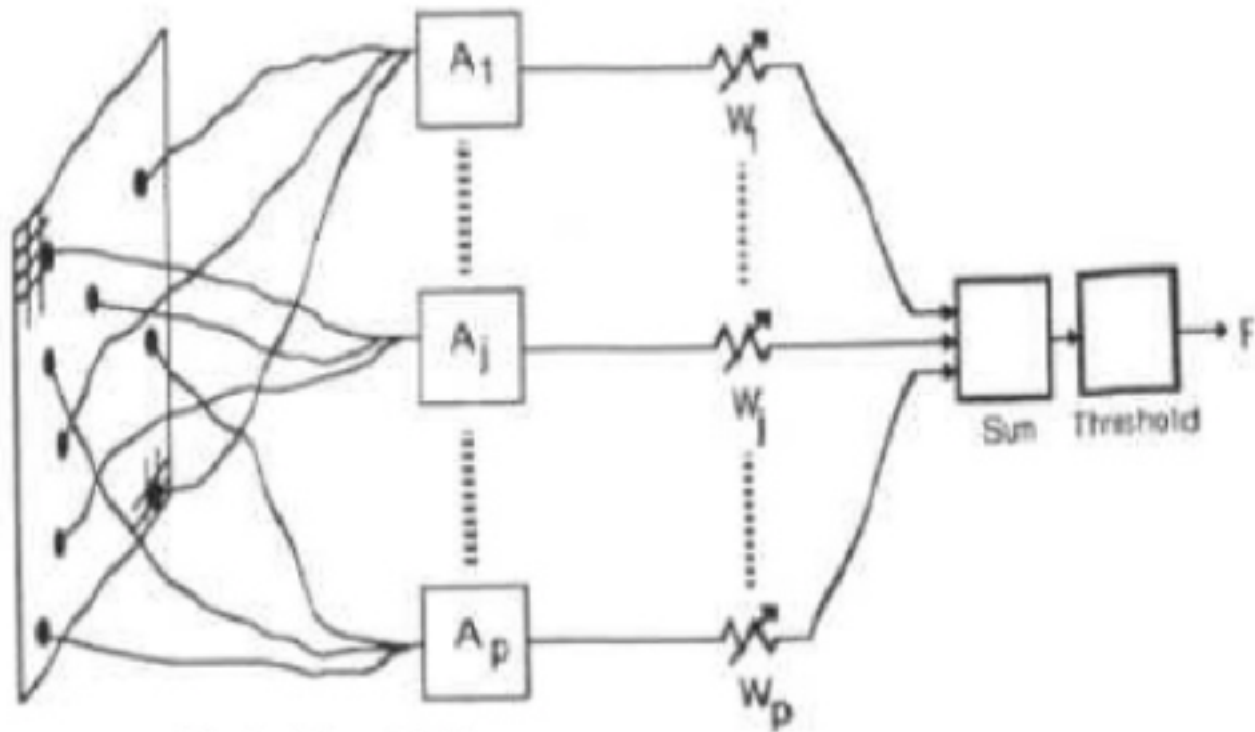
sigmoid



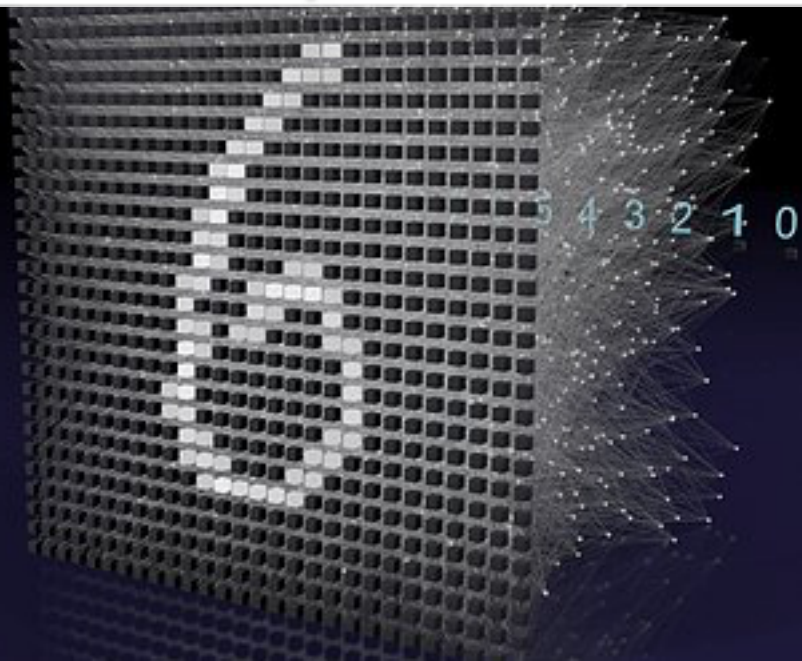
$$a = \frac{1}{1 + \exp(-z)}$$



# The Mark 1



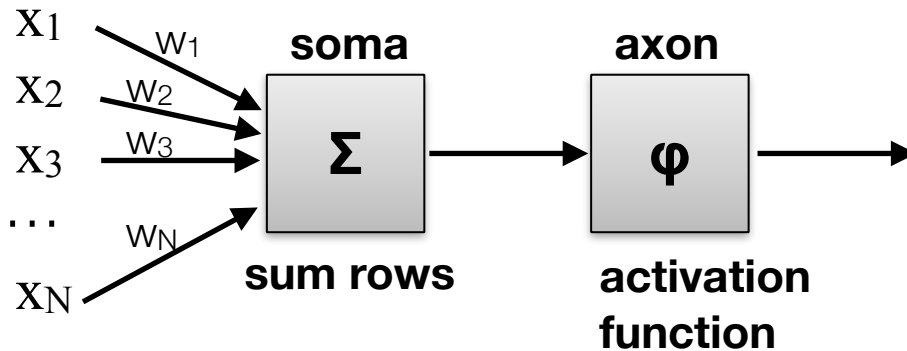
PERCEPTRON



Perceptron Learning Rule:  
~Stochastic Gradient Descent

# Layers Notation

dendrite



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

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

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

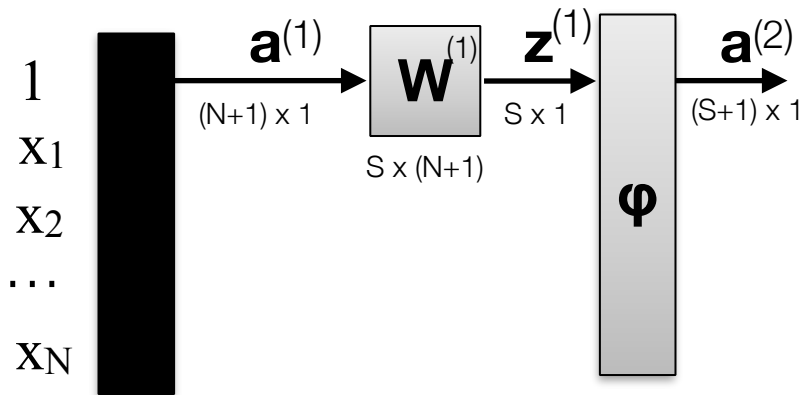
$$\mathbf{W} = \begin{bmatrix} w_{0,1} & w_{0,2} & \dots & w_{0,N} \\ w_{1,1} & w_{1,2} & \dots & w_{1,N} \\ \vdots & & & \\ w_{S,1} & w_{S,2} & \dots & w_{S,N} \end{bmatrix}$$

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

$\mathbf{a}^{next} = \phi(\mathbf{z}^{current})$ , concat bias term

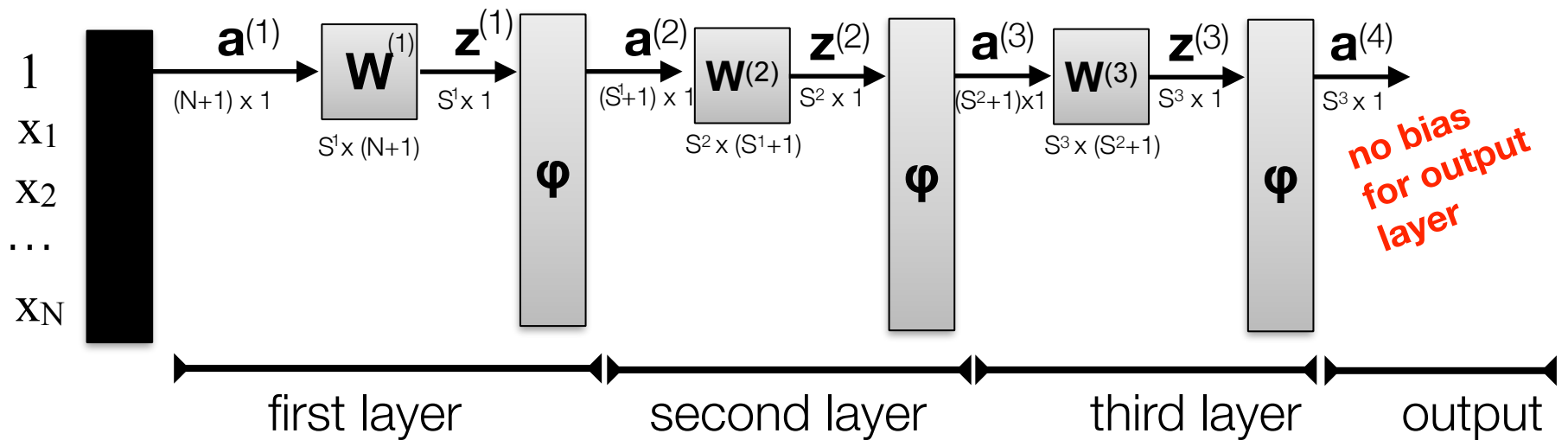
$$\mathbf{a}^{(next)} = \begin{bmatrix} 1 \\ \phi(z_1^{curr}) \\ \vdots \\ \phi(z_N^{curr}) \end{bmatrix} \rightarrow \mathbf{a}^{(L)} = \begin{bmatrix} 1 \\ \phi(z_1^{L-1}) \\ \vdots \\ \phi(z_N^{L-1}) \end{bmatrix}$$

input



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

# Generic Multiple Layers Notation



$\mathbf{a}^{(L+1)} = \phi(\mathbf{z}^{(L)})$ , concat bias term       $\mathbf{a}^{(final)}$  rows=unique classes

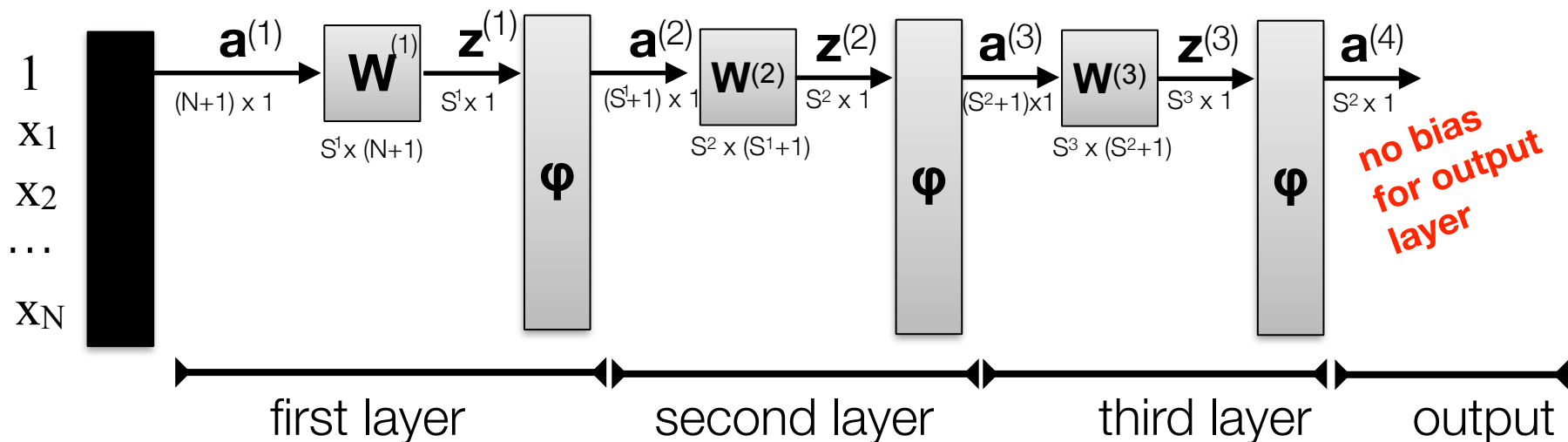
$$\mathbf{z}^{(L)} = \mathbf{W}^{(L)} \cdot \mathbf{a}^{(L)}$$

$$\mathbf{W}^{(L)} = \begin{bmatrix} w_{0,1}^{(L)} & w_{0,2}^{(L)} & \dots & w_{0,S^L}^{(L)} \\ w_{1,1}^{(L)} & w_{1,2}^{(L)} & \dots & w_{1,S^L}^{(L)} \\ \vdots & & & \\ w_{S^{L+1},1}^{(L)} & w_{S^{L+1},2}^{(L)} & \dots & w_{S^{L+1},S^L}^{(L)} \end{bmatrix}$$

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

$$\mathbf{z}^{(L)} = \mathbf{W}^{(L)} \cdot \phi(\mathbf{W}^{(L-1)} \cdot \phi(\mathbf{z}^{(L-2)}))$$

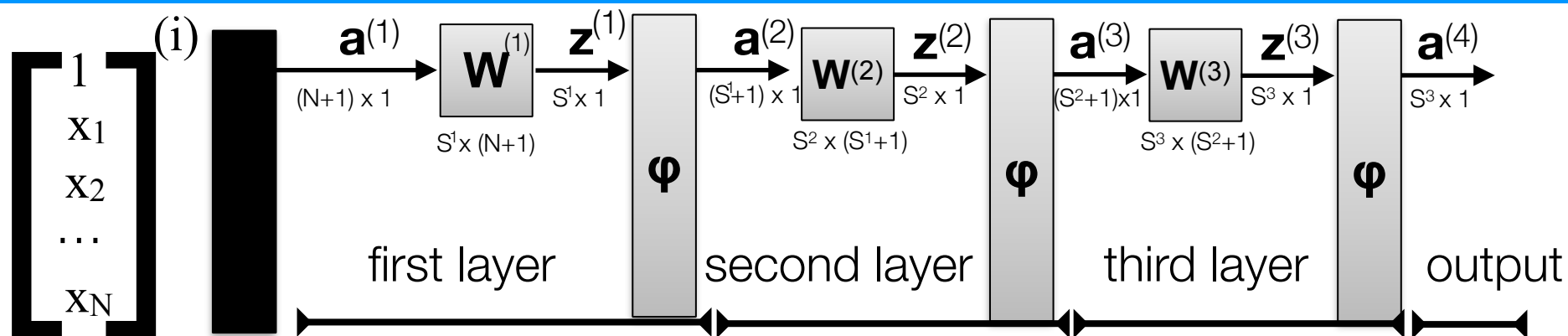
# Multiple layers notation



- **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 12

# Compact feedforward notation



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

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

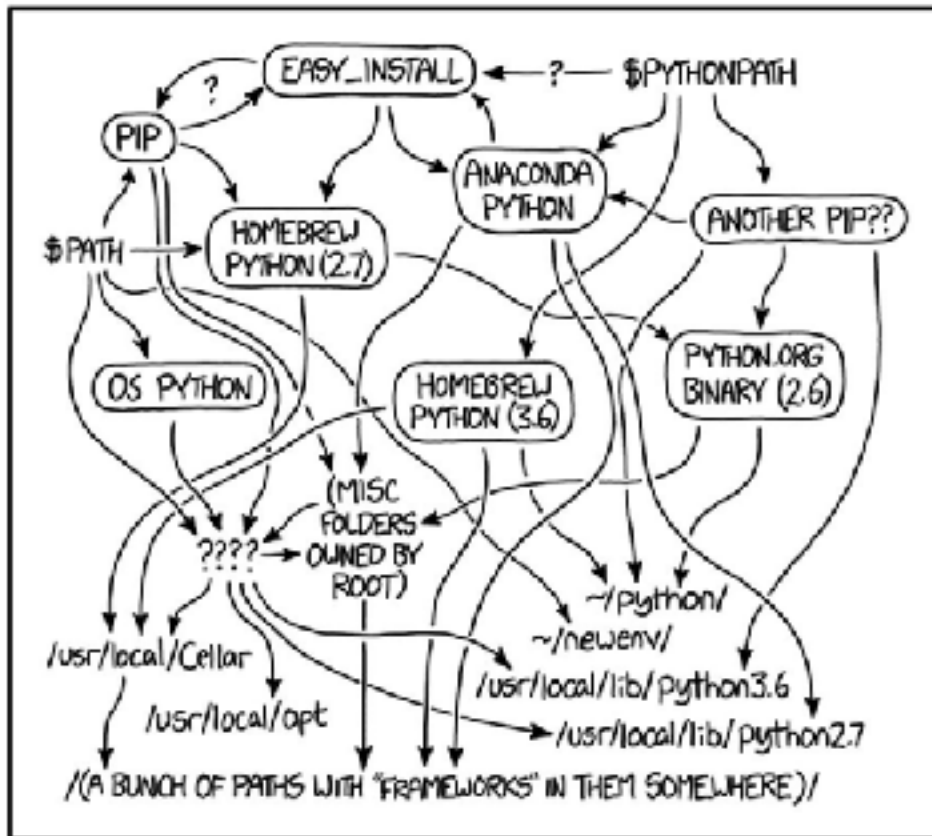
$$\begin{bmatrix} z^{(L)}_1 \\ \vdots \\ z^{(L)}_{S^L} \end{bmatrix}^{(i)} = [\mathbf{W}^{(L)}] \begin{bmatrix} a^{(L)}_1 \\ \vdots \\ a^{(L)}_{S^{L-1}+1} \end{bmatrix}^{(i)}$$

$$\begin{bmatrix} \begin{bmatrix} z^{(L)}_1 \\ \vdots \\ z^{(L)}_{S^L} \end{bmatrix}^{(1)} & \dots & \begin{bmatrix} z^{(L)}_1 \\ \vdots \\ z^{(L)}_{S^L} \end{bmatrix}^{(M)} \end{bmatrix} = [\mathbf{W}^{(L)}] \begin{bmatrix} \begin{bmatrix} a^{(L)}_1 \\ \vdots \\ a^{(L)}_{S^{L-1}+1} \end{bmatrix}^{(1)} & \dots & \begin{bmatrix} a^{(L)}_1 \\ \vdots \\ a^{(L)}_{S^{L-1}+1} \end{bmatrix}^{(M)} \end{bmatrix}$$

$$\mathbf{Z}^{(L)} = \mathbf{W}^{(L)} \cdot \mathbf{A}^{(L)}$$

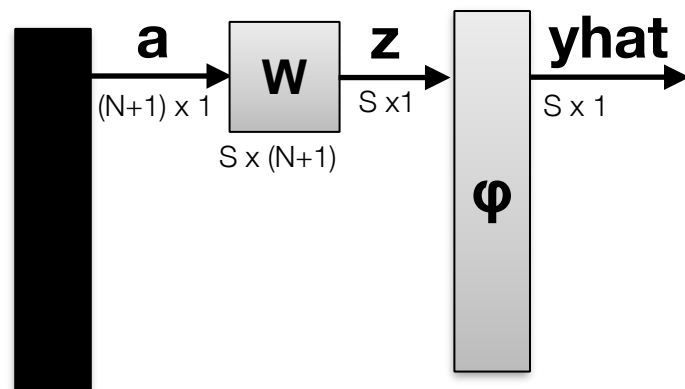
$$\mathbf{Z}^{(L)} = \mathbf{W}^{(L)} \cdot \phi(\mathbf{Z}^{(L-1)})$$

# Training Neural Network Architectures



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

# Rosenblatt's Perceptron, 1957



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



Need objective Function, minimize MSE  $J(\mathbf{W}) = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$

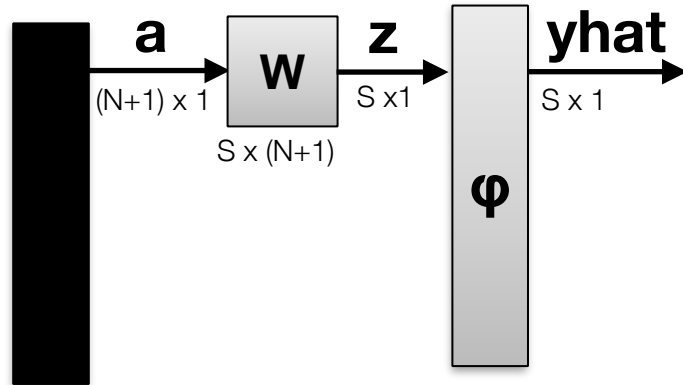
where ground truth  $\mathbf{y}^{(i)}$  is one-hot encoded!

$$\begin{bmatrix} y_1 \\ \dots \\ y_c \end{bmatrix}^{(i)} \rightarrow \begin{bmatrix} \begin{bmatrix} y_1 \\ \dots \\ y_c \end{bmatrix}^{(1)} \quad \dots \quad \begin{bmatrix} y_1 \\ \dots \\ y_c \end{bmatrix}^{(M)} \end{bmatrix} = \mathbf{Y}$$



# Simple Architectures

- Rosenblatt's perceptron, 1957



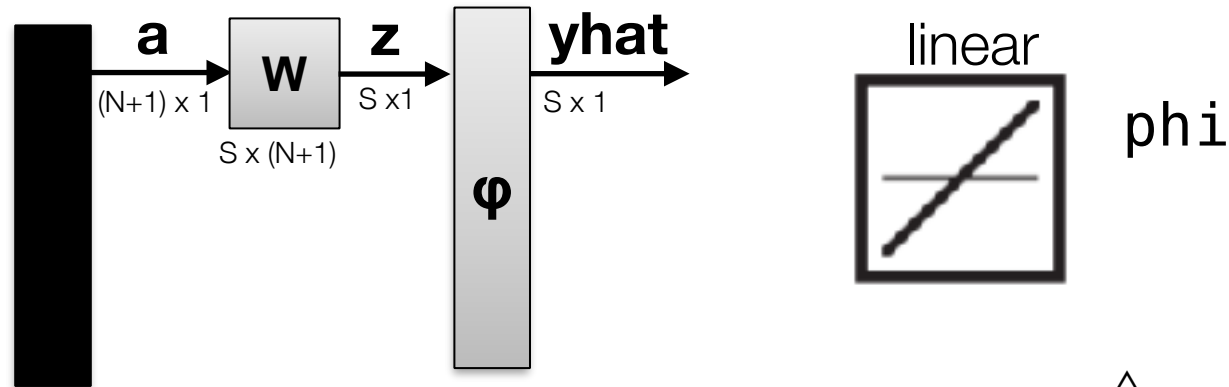
**Self Test** - If this is a binary classification problem, how large is  $S$ , the length of  $\mathbf{\hat{y}}$ ?

- A. Can't determine
- B. 2
- C. 1
- D.  $N$

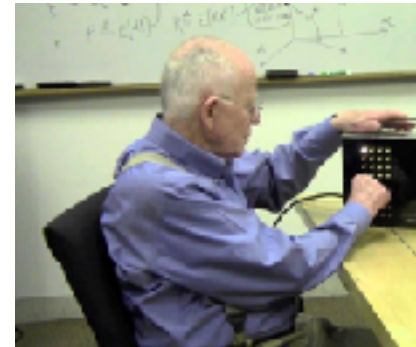


# Simple Architectures

- Adaline network, Widrow and Hoff, 1960



Marcian "Ted" Hoff



Bernard Widrow

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

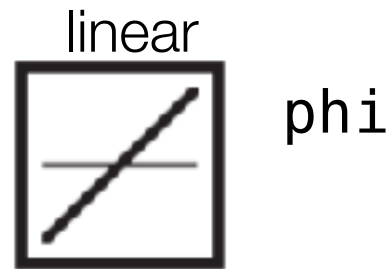
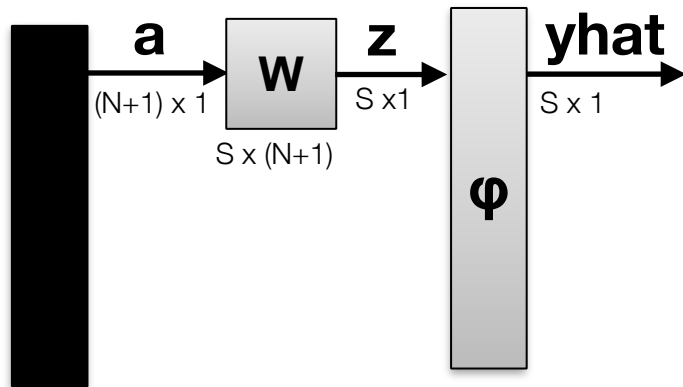
New objective function becomes:  $J(\mathbf{W}) = ||\mathbf{Y} - \mathbf{W} \cdot \mathbf{A}||^2$

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

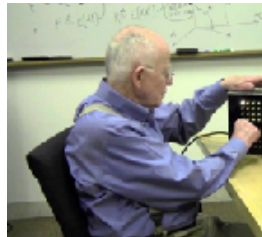
We have been using the **Widrow-Hoff Learning Rule**

# Simple Architectures

- Adaline network, Widrow and Hoff, 1960



Marcian "Ted" Hoff



Bernard Widrow

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

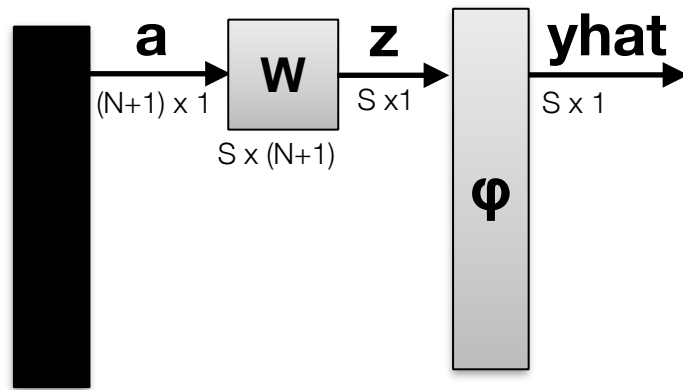
For case  $S=1$ ,  $\mathbf{W}$  has only one row,  $\mathbf{w}$   
this is just **linear regression**...

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



# Simple Architectures

- Modern Perceptron network



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

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

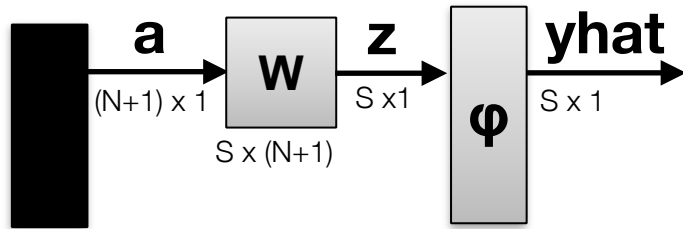
For case  $S=1$ , this is just **logistic regression...**  
and **we have already solved this!**

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



What happens when  $S > 1$  ?

# What if we have more than $S=1$ ?



$$\begin{bmatrix} \begin{bmatrix} \phi(z_1) \\ \dots \\ \phi(z_s) \end{bmatrix}^{(1)} & \dots & \begin{bmatrix} \phi(z_1) \\ \dots \\ \phi(z_s) \end{bmatrix}^{(M)} \end{bmatrix} = \hat{\mathbf{Y}}^{\Delta}$$

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

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

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

$$\mathbf{yhat}^{(i)} = \begin{bmatrix} \phi(\text{row}=1 \mathbf{W} \cdot \mathbf{x}^{(i)}) \\ \phi(\text{row}=2 \mathbf{W} \cdot \mathbf{x}^{(i)}) \\ \dots \\ \phi(\text{row}=S \mathbf{W} \cdot \mathbf{x}^{(i)}) \end{bmatrix}$$

one hot



which is one-versus-all!

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

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

...

$$J(s\mathbf{W}) = \sum_{i=1} [y_s^{(i)} - \phi(s\mathbf{W} \cdot \mathbf{x}^{(i)})]^2$$

# 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  $\mathbf{w}_{\text{class}}$  be rows of weight matrix,  $\mathbf{W}$ 
  - 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  $\alpha$ -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

TED2008 · 13:33 · Filmed Feb 2008

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

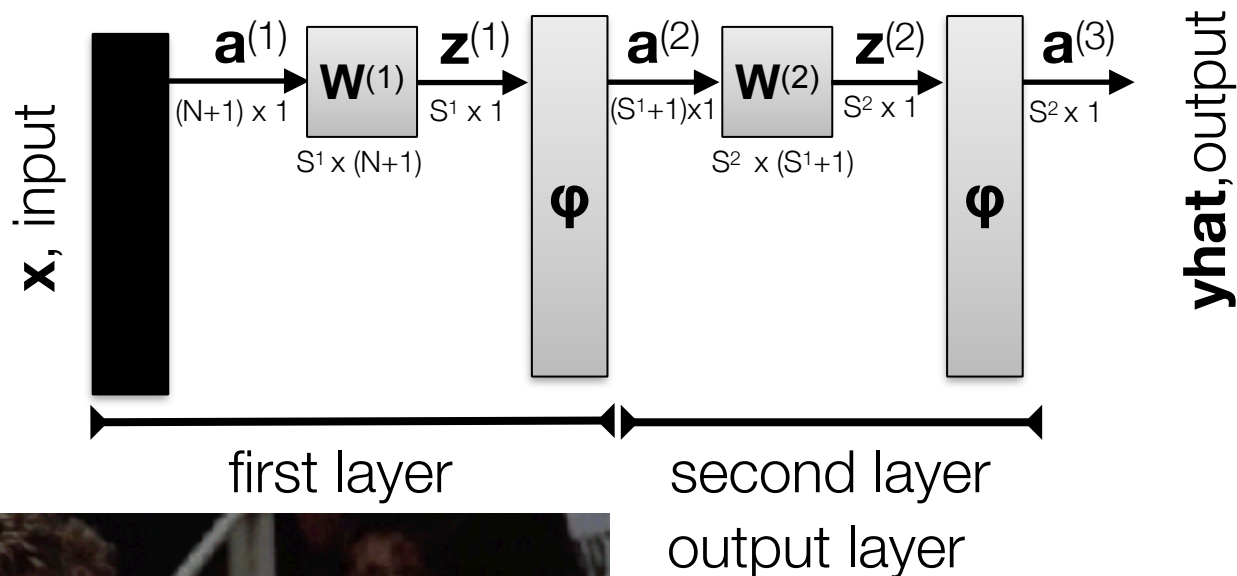


Geoffrey Hinton



# 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!!!



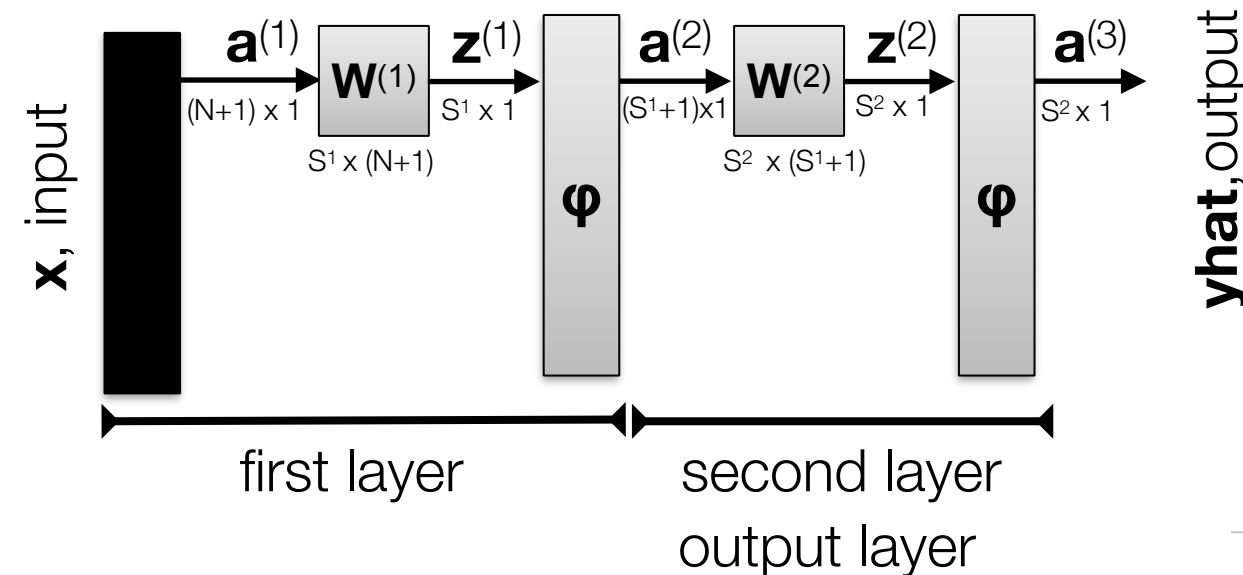
$$\mathbf{yhat}^{(i)} = \begin{bmatrix} \phi(\text{row}=1 \mathbf{W}^{(2)} \cdot \phi(\mathbf{W}^{(1)} \mathbf{a}^{(1)})) \\ \vdots \\ \phi(\text{row}=S \mathbf{W}^{(2)} \cdot \phi(\mathbf{W}^{(1)} \mathbf{a}^{(1)})) \end{bmatrix}$$

one hot



# Back propagation

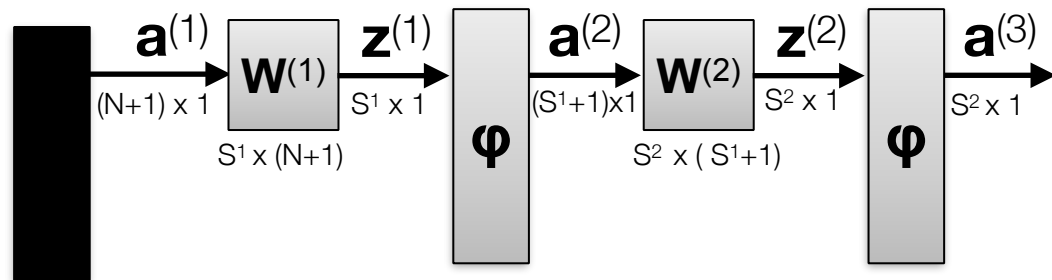
- Steps:
  - propagate weights forward
  - calculate gradient at final layer
  - back propagate gradient for each layer
    - via recurrence relation



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

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

# Back propagation



use chain rule:

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

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

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

Solve this in explainer video for next in class assignment!

# End of Session

- 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=0ahUKEwioprwn27fPAhWMx4MKHYbwDIwQFggeMAA&url=http%3A%2F%2Fhagan.okstate.edu%2FNNDesign.pdf&usg=AFQjCNG5YbM4xSMm6K5HNsG-4Q8TvOu\\_Lw&sig2=bgT3k-5ZDDTPZ07Qu8Oreg](https://www.google.com/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=0ahUKEwioprwn27fPAhWMx4MKHYbwDIwQFggeMAA&url=http%3A%2F%2Fhagan.okstate.edu%2FNNDesign.pdf&usg=AFQjCNG5YbM4xSMm6K5HNsG-4Q8TvOu_Lw&sig2=bgT3k-5ZDDTPZ07Qu8Oreg)

Michael Nielsen

<http://neuralnetworksanddeeplearning.com>