

Bil 470 / YAP 470

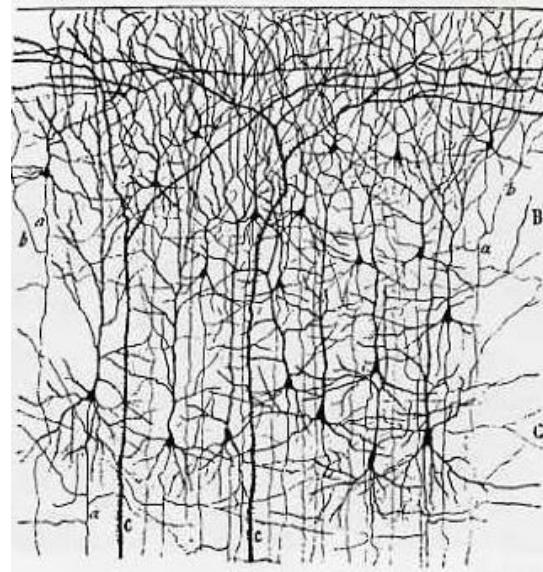
Introduction to Machine Learning (Yapay Öğrenme)

Batuhan Bardak

Neural Networks

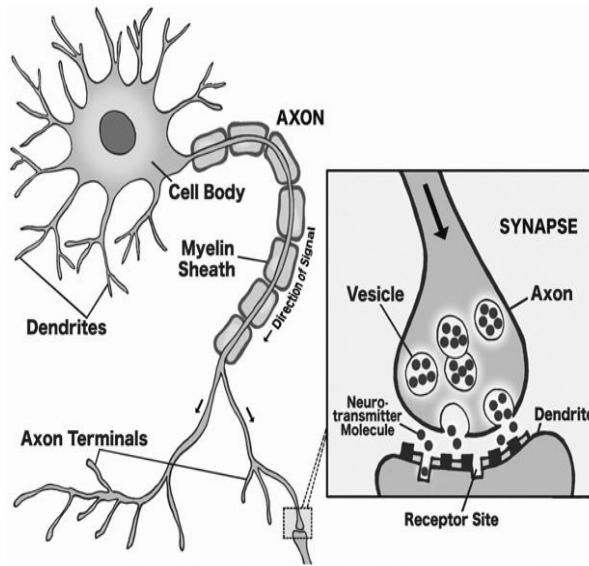
Date: 10.10.2022

Discovery of Neurons



Ramón y Cajal, 1900

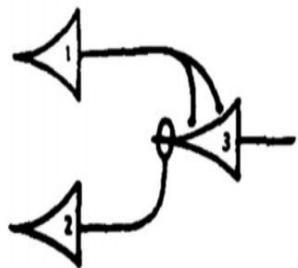
How the Human Brain Learns



Inputs arrive at dendrites, and axon serve as output channel

A Mathematical Model of Brain: McCulloch & Pitts Neuron

Artificial Neurons



$$N_1(t) \cdot \equiv \cdot N_1(t - 1) \cdot \sim N_2(t - 1)$$

MCP, 1943

BULLETIN OF
MATHEMATICAL BIOPHYSICS
VOLUME 5, 1943

A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY

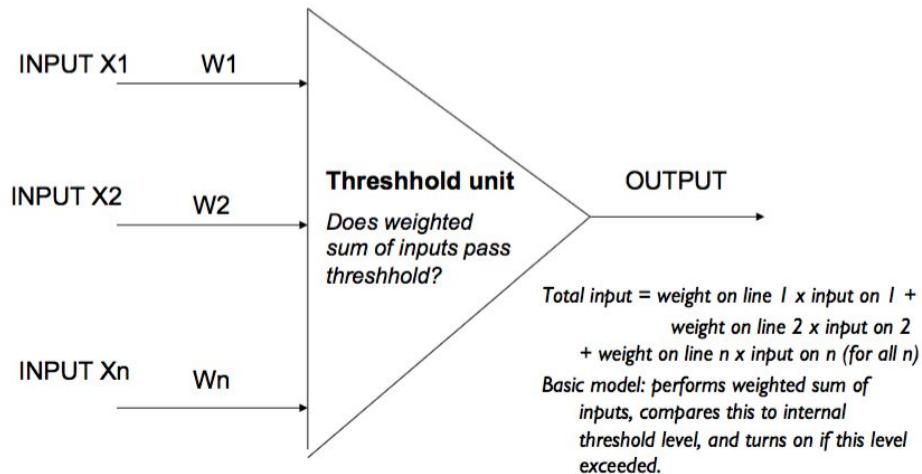
WARREN S. MCCULLOCH AND WALTER PITTS

FROM THE UNIVERSITY OF ILLINOIS, COLLEGE OF MEDICINE,
DEPARTMENT OF PSYCHIATRY AT THE ILLINOIS NEUROPSYCHIATRIC INSTITUTE,
AND THE UNIVERSITY OF CHICAGO

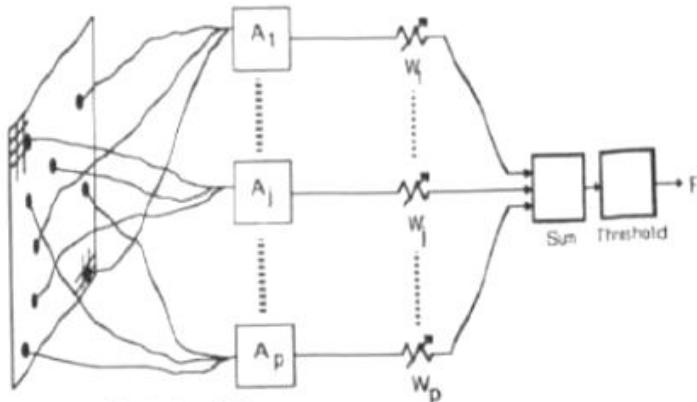
Because of the "all-or-none" character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every net can be described in these terms, with the addition of more complicated logical means for nets containing circles; and that for any logical expression satisfying certain conditions, one can find a net behaving in the fashion it describes. It is shown that many particular choices among possible neurophysiological assumptions are equivalent, in the sense that for every net behaving under one assumption, there exists another net which behaves under the other and gives the same results, although perhaps not in the same time. Various applications of the calculus are discussed.

McCulloch & Pitts Neuron

McCulloch Pitts Neuron



Perceptron

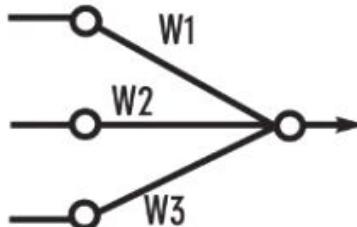


Frank Rosenblatt
(1928-1971)

Original Perceptron

*(From Perceptrons by M. L Minsky and S. Papert,
1969, Cambridge, MA: MIT Press. Copyright 1969
by MIT Press.)*

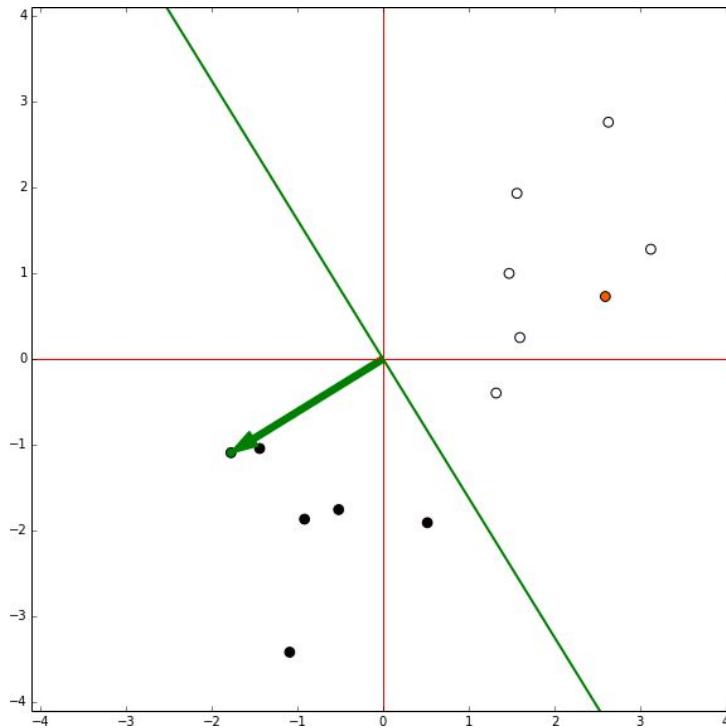
Simplified model:



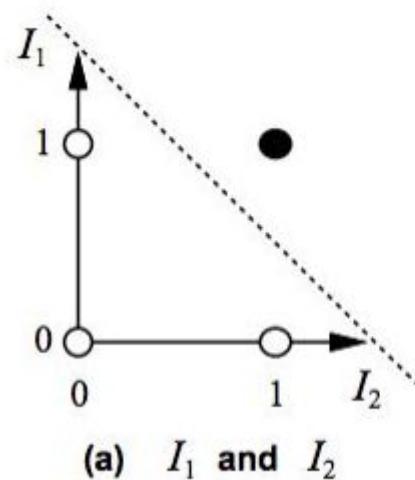
23

Frank Rosenblatt (1957)

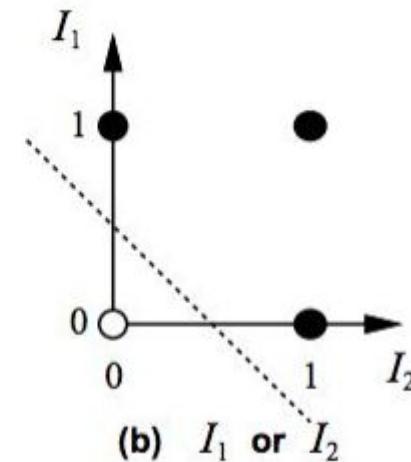
Perceptron Learning Rule



Linear Seperable Case



(a) I_1 and I_2

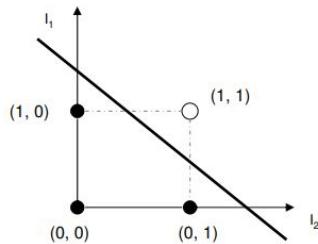


(b) I_1 or I_2

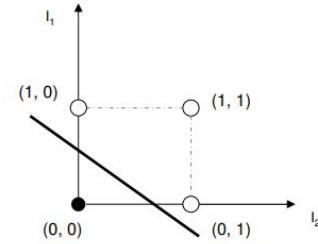
Xor Case

Use more than one neuron!

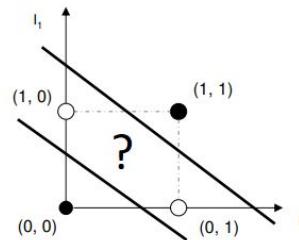
AND		
I_1	I_2	out
0	0	0
0	1	0
1	0	0
1	1	1



OR		
I_1	I_2	out
0	0	0
0	1	1
1	0	1
1	1	1

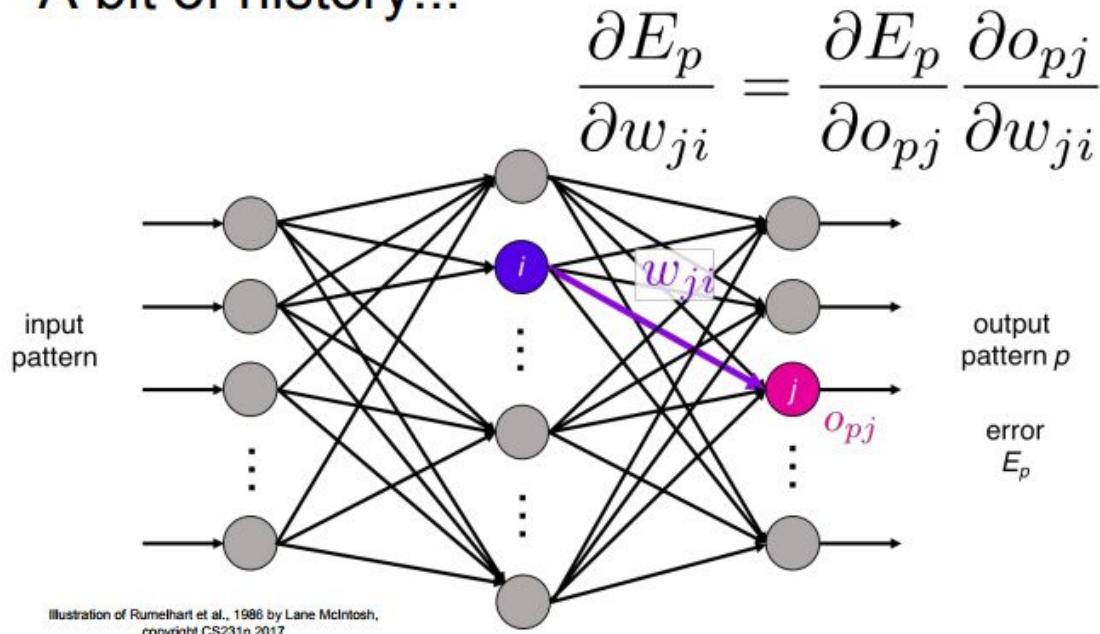


XOR		
I_1	I_2	out
0	0	0
0	1	1
1	0	1
1	1	0



Backpropagation (1985)

A bit of history...



$$\frac{\partial E_p}{\partial w_{ji}} = \frac{\partial E_p}{\partial o_{pj}} \frac{\partial o_{pj}}{\partial w_{ji}}$$

nature > letters > article

nature
International journal of science

MENU ▾ Letter | Published: 09 October 1986

Learning representations by back-propagating errors

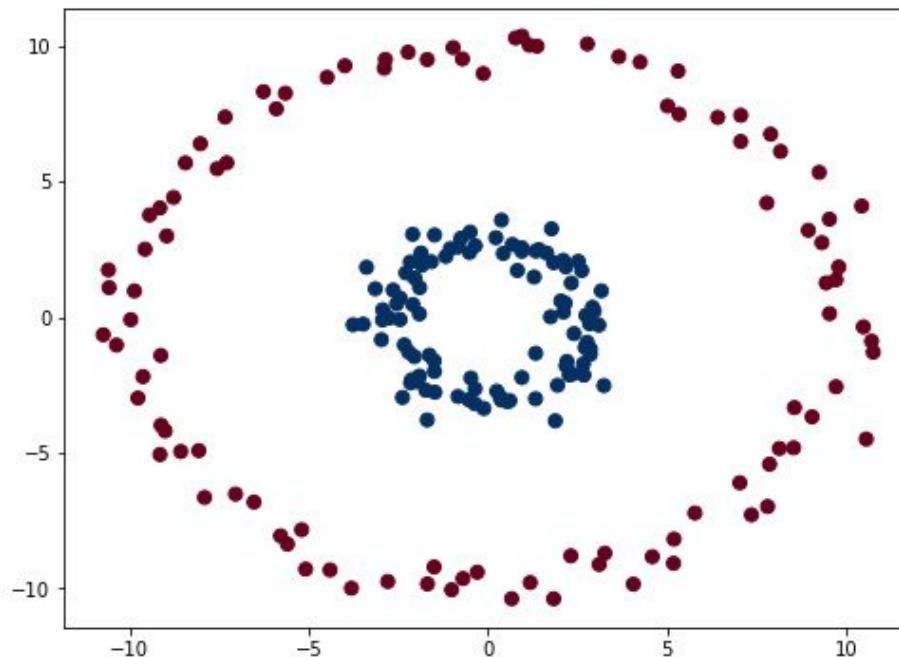
David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

Nature 323, 533–536 (09 October 1986) | Download Citation ↴

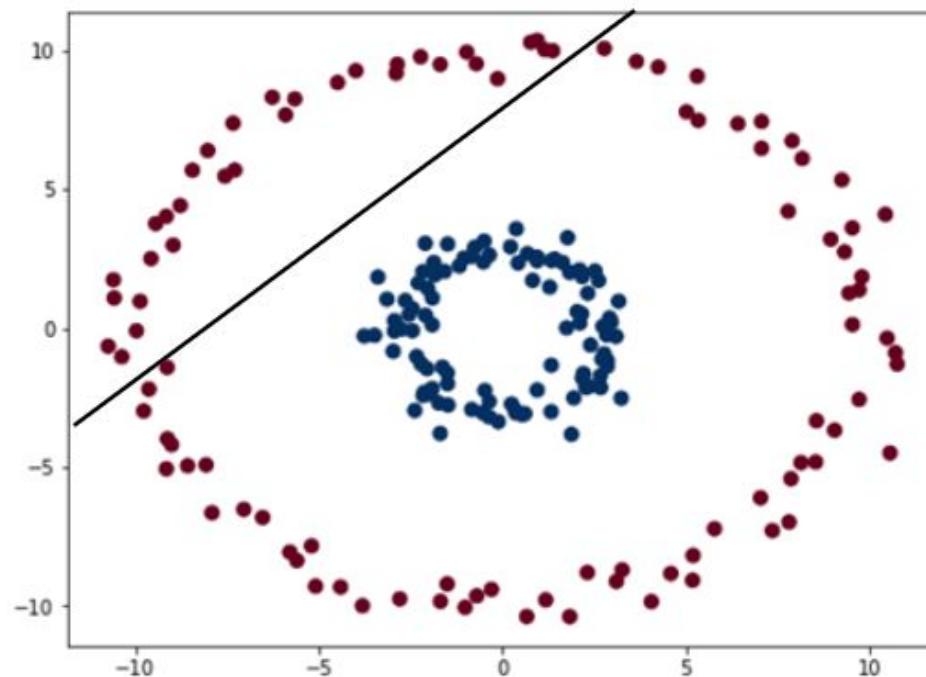
Abstract

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

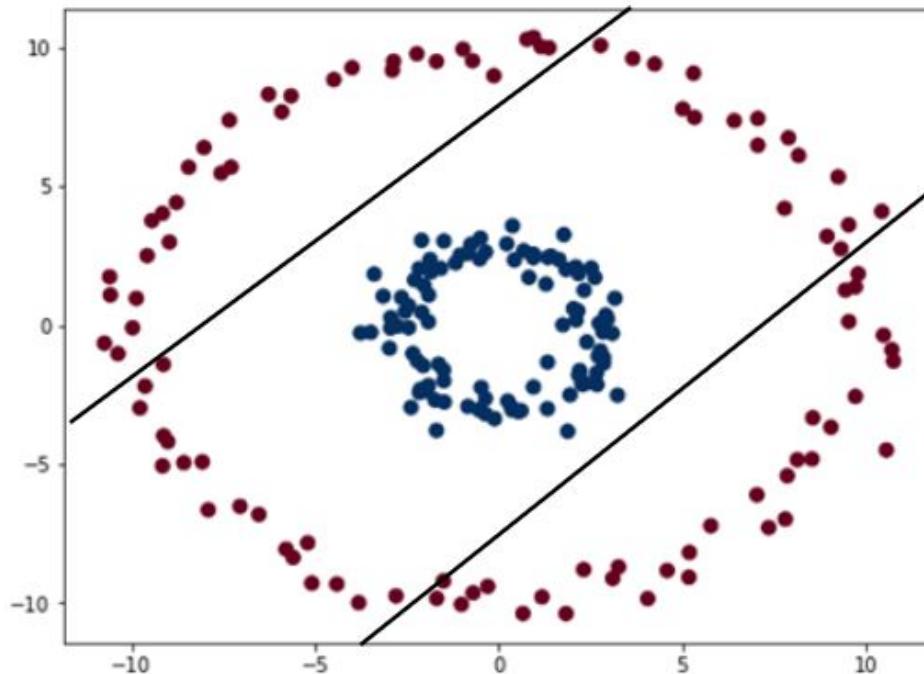
What to do in a More Complex Scenario?



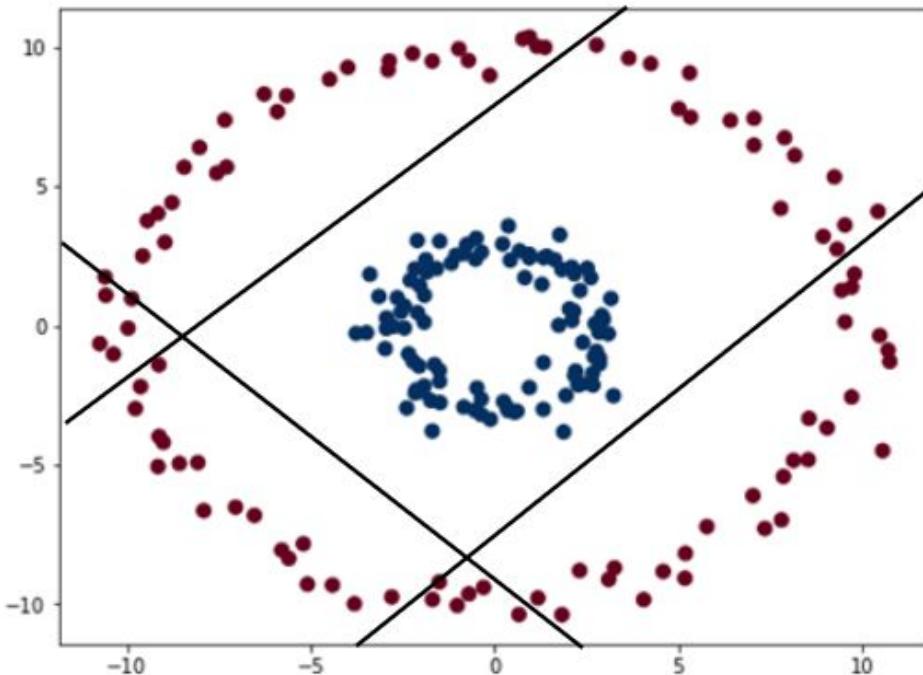
First Neuron



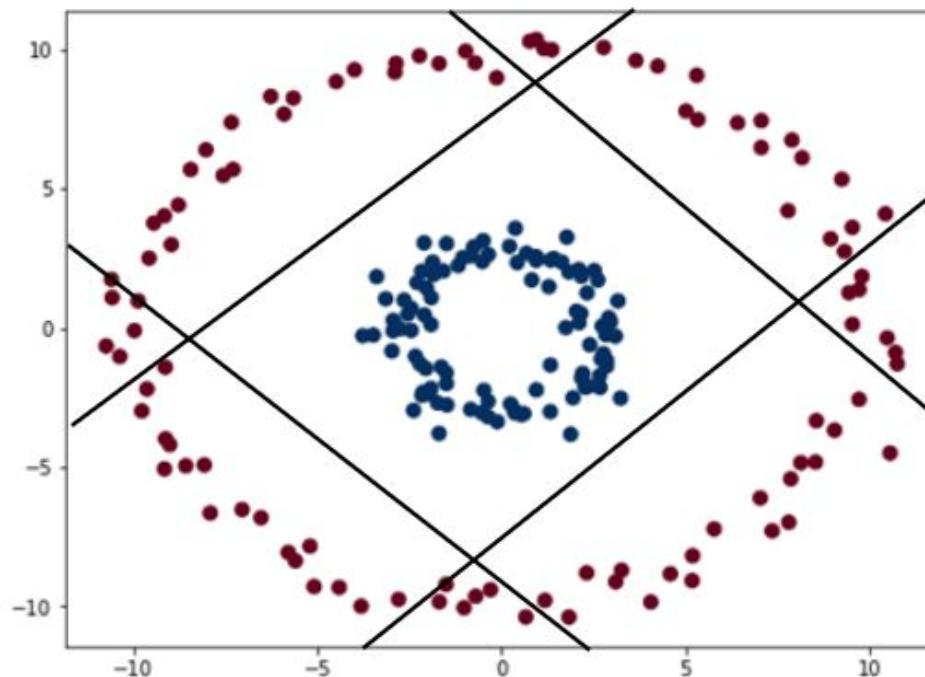
Second Neuron



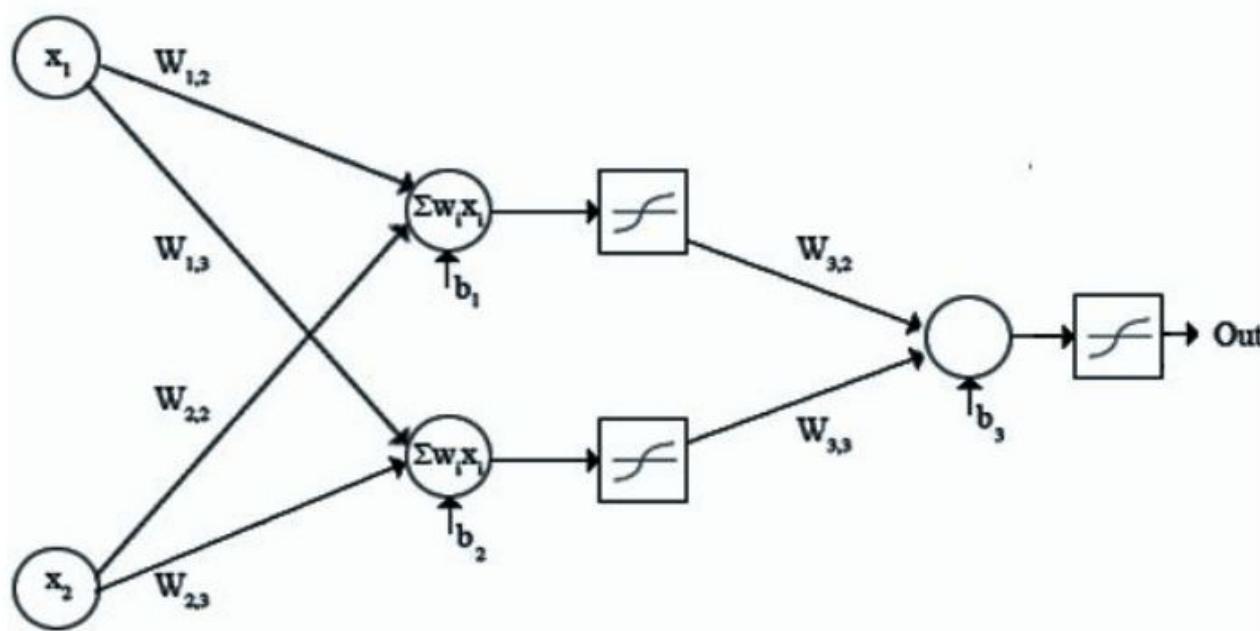
Third Neuron



Fourth Neuron

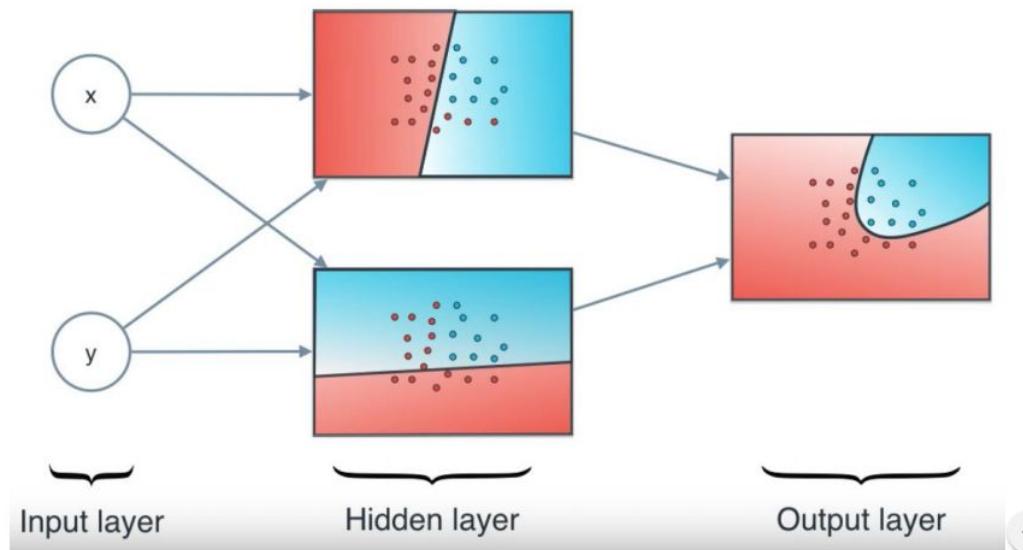


Multilayer Perceptron



Multilayer Perceptron

Multilayer Perceptron



Neural Network for classification

Vector function with tunable parameters θ

$$\mathbf{f}(\cdot; \theta) : \mathbb{R}^N \rightarrow (0, 1)^K$$

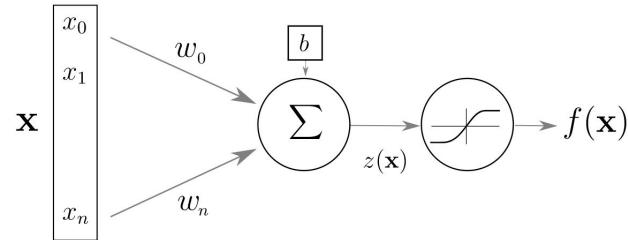
Sample s in dataset S :

- input: $\mathbf{x}^s \in \mathbb{R}^N$
- expected output: $y^s \in [0, K - 1]$

Output is a conditional probability distribution:

$$\mathbf{f}(\mathbf{x}^s; \theta)_c = P(Y = c | X = \mathbf{x}^s)$$

Artificial Neuron

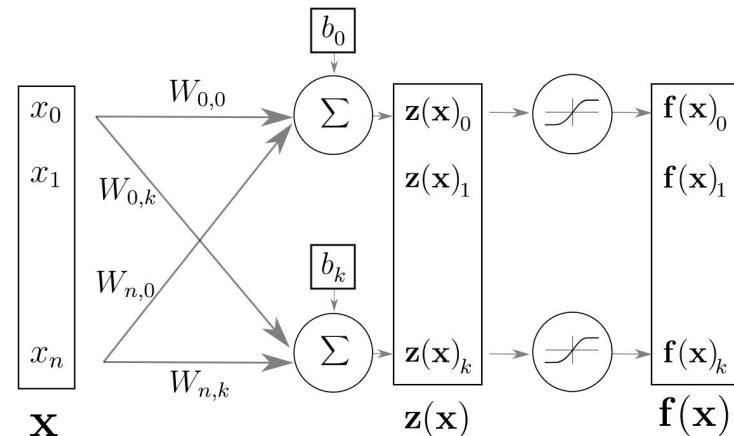


$$z(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

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

- $\mathbf{x}, f(\mathbf{x})$ input and output
- $z(\mathbf{x})$ pre-activation
- \mathbf{w}, b weights and bias
- g activation function

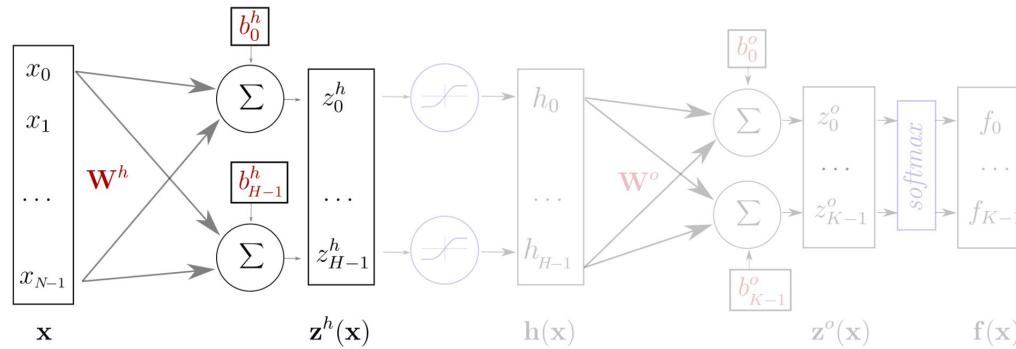
Layer of Neurons



$$\mathbf{f}(\mathbf{x}) = g(\mathbf{z}(\mathbf{x})) = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

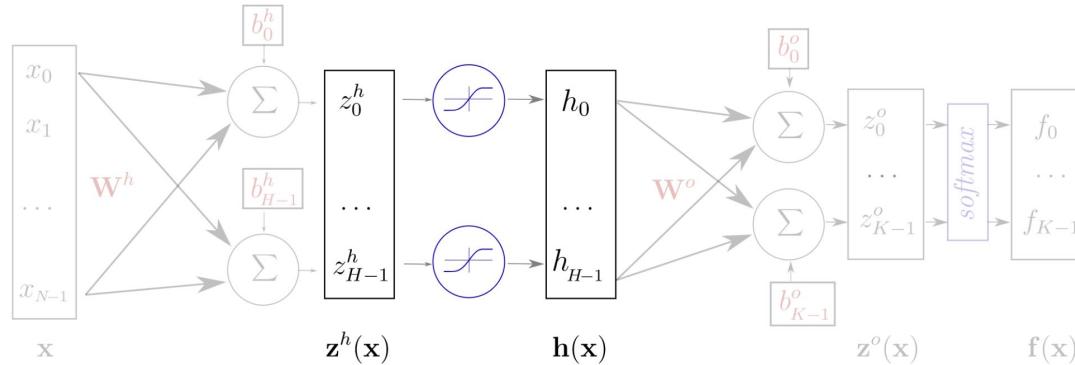
- \mathbf{W}, \mathbf{b} now matrix and vector

One Hidden Layer Network



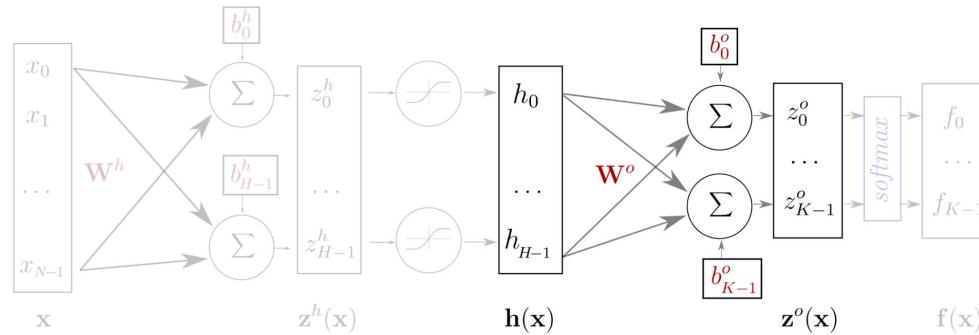
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = softmax(\mathbf{z}^o) = softmax(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

One Hidden Layer Network



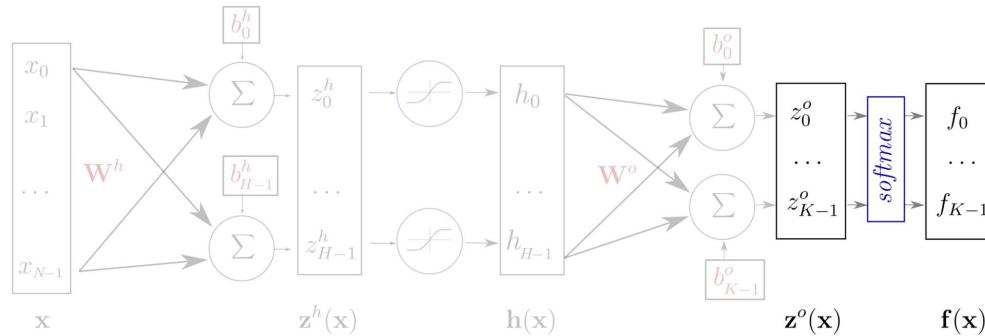
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = softmax(\mathbf{z}^o) = softmax(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

One Hidden Layer Network



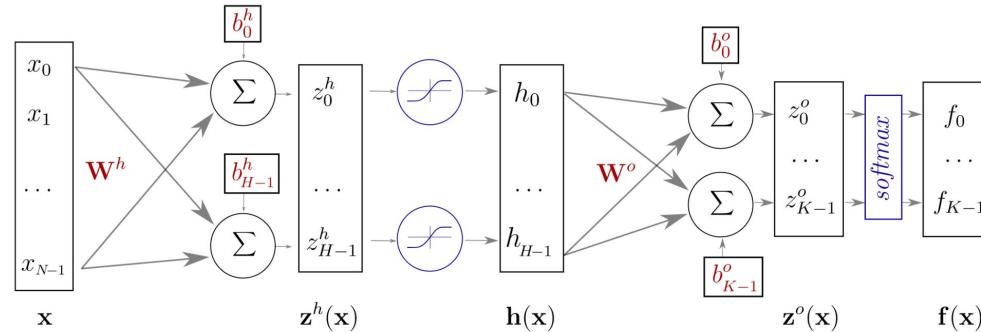
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = softmax(\mathbf{z}^o) = softmax(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

One Hidden Layer Network

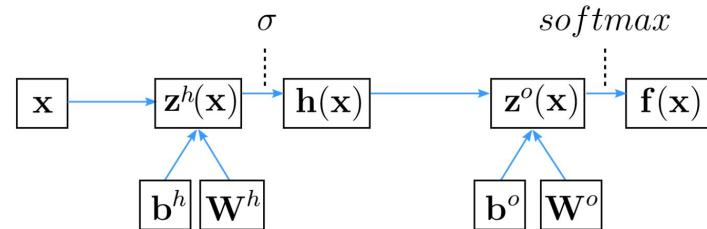


- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

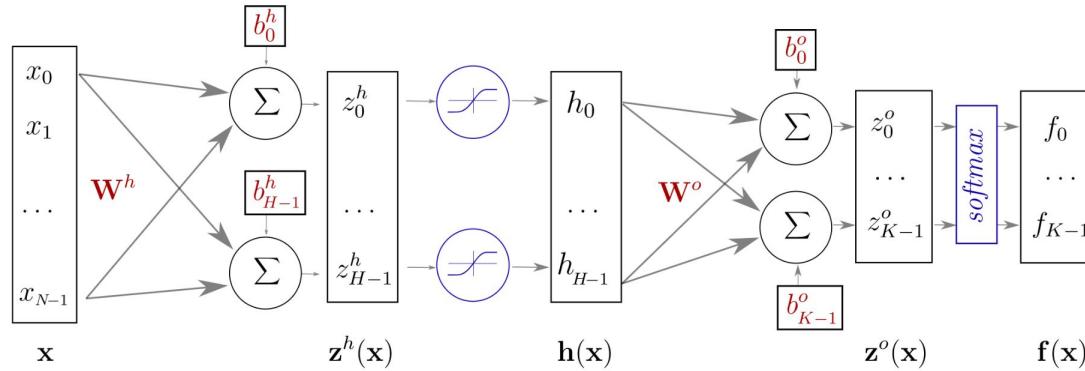
One Hidden Layer Network



Alternate representation



One Hidden Layer Network



Keras implementation

```
model = Sequential()
model.add(Dense(H, input_dim=N)) # weight matrix dim [N * H]
model.add(Activation("tanh"))
model.add(Dense(K))           # weight matrix dim [H x K]
model.add(Activation("softmax"))
```

Softmax Function

$$\text{softmax}(\mathbf{x}) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

$$\frac{\partial \text{softmax}(\mathbf{x})_i}{\partial x_j} = \begin{cases} \text{softmax}(\mathbf{x})_i \cdot (1 - \text{softmax}(\mathbf{x})_i) & i = j \\ -\text{softmax}(\mathbf{x})_i \cdot \text{softmax}(\mathbf{x})_j & i \neq j \end{cases}$$

- vector of values in $(0, 1)$ that add up to 1
- $p(Y = c | X = \mathbf{x}) = \text{softmax}(\mathbf{z}(\mathbf{x}))_c$
- the pre-activation vector $\mathbf{z}(\mathbf{x})$ is often called "the logits"

Training the Network

Find parameters $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$ that minimize the
negative log likelihood (or [cross entropy](#))

The loss function for a given sample $s \in S$:

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = nll(\mathbf{x}^s, y^s; \theta) = -\log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

example $y^s = 3$

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = l \left(\begin{pmatrix} f_0 \\ \vdots \\ f_3 \\ \vdots \\ f_{K-1} \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \right) = -\log f_3$$

Training the Network

Find parameters $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$ that minimize the **negative log likelihood** (or [cross entropy](#))

The loss function for a given sample $s \in S$:

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = nll(\mathbf{x}^s, y^s; \theta) = -\log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

The cost function is the negative likelihood of the model computed on the full training set (for i.i.d. samples):

$$L_S(\theta) = -\frac{1}{|S|} \sum_{s \in S} \log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s} + \lambda \Omega(\theta)$$

$\lambda \Omega(\theta) = \lambda(||W^h||^2 + ||W^o||^2)$ is an optional regularization term.

Stochastic Gradient Descent

Initialize θ randomly

For E epochs perform:

- Randomly select a small batch of samples ($B \subset S$)
 - Compute gradients: $\Delta = \nabla_{\theta} L_B(\theta)$
 - Update parameters: $\theta \leftarrow \theta - \eta \Delta$
 - $\eta > 0$ is called the learning rate
- Repeat until the epoch is completed (all of S is covered)

Stop when reaching criterion:

- nll stops decreasing when computed on validation set

Computing Gradients

Output Weights: $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial W_{i,j}^o}$

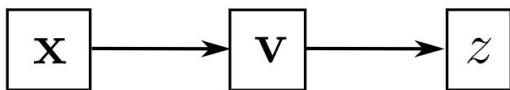
Hidden Weights: $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial W_{i,j}^h}$

Output bias: $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial b_i^o}$

Hidden bias: $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial b_i^h}$

- The network is a composition of differentiable modules
- We can apply the "chain rule"

Chain rule



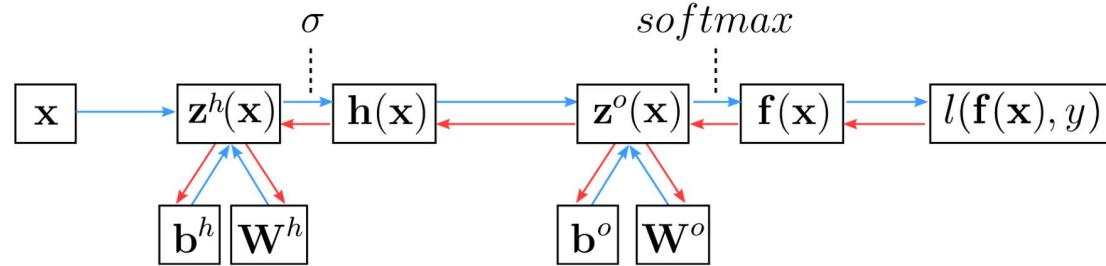
$$z = u(\mathbf{v}(\mathbf{x}))$$

chain-rule

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial v_j} \frac{\partial v_j}{\partial x_i} = \nabla u \cdot \frac{\partial \mathbf{v}}{\partial x_i}$$

$$\begin{array}{c|c|c} v_j & \frac{\partial v_j}{\partial x_i} & \frac{\partial z}{\partial v_j} \\ \hline \mathbf{v} & \frac{\partial \mathbf{v}}{\partial x_i} & \nabla u \end{array}$$

Backpropagation



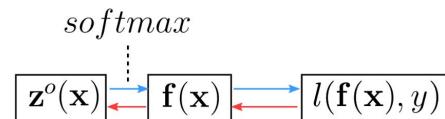
Compute partial derivatives of the loss

$$\bullet \frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{\partial -\log \mathbf{f}(\mathbf{x})_y}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{-1_{y=i}}{\mathbf{f}(\mathbf{x})_y} = \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_i}$$

$$\bullet \frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} = ?$$

Backpropagation

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} &= \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i} & \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_i} &= \frac{-1_{y=i}}{\mathbf{f}(\mathbf{x})_y} \\ &= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial softmax(\mathbf{z}^o(\mathbf{x}))_j}{\partial \mathbf{z}^o(\mathbf{x})_i} & \mathbf{f}(\mathbf{x}) &= softmax(\mathbf{z}^o(\mathbf{x}))\end{aligned}$$

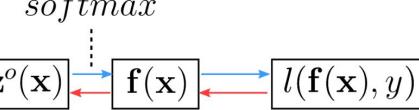


Backpropagation

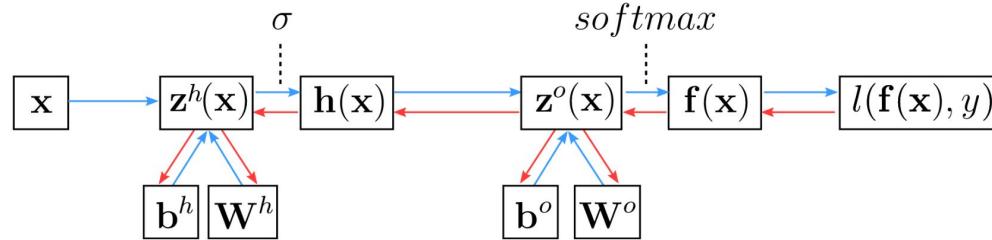
$$\begin{aligned}
 \frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} &= \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
 &= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
 &= -\frac{1}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_y}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
 &= \begin{cases} -\frac{1}{\mathbf{f}(\mathbf{x})_y} \text{softmax}(\mathbf{z}^o(\mathbf{x}))_y (1 - \text{softmax}(\mathbf{z}^o(\mathbf{x}))_y) & \text{if } i = y \\ \frac{1}{\mathbf{f}(\mathbf{x})_y} \text{softmax}(\mathbf{z}^o(\mathbf{x}))_y \text{softmax}(\mathbf{z}^o(\mathbf{x}))_i & \text{if } i \neq y \end{cases} \\
 &= \begin{cases} -1 + \mathbf{f}(\mathbf{x})_y & \text{if } i = y \\ \mathbf{f}(\mathbf{x})_i & \text{if } i \neq y \end{cases}
 \end{aligned}$$

$$\nabla_{\mathbf{z}^o(\mathbf{x})} l(\mathbf{f}(\mathbf{x}), y) = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$$

$\mathbf{e}(y)$: one-hot encoding of y



Backpropagation



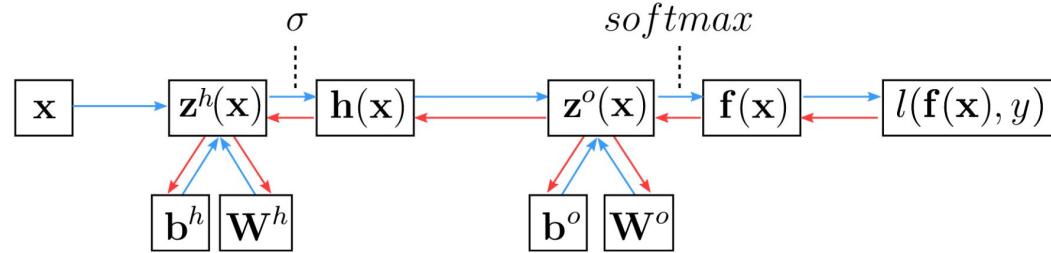
Gradients

- $\nabla_{\mathbf{z}^o(\mathbf{x})} l = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

- $\nabla_{\mathbf{b}^o} l = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

because $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$ and then $\frac{\partial \mathbf{z}^o(\mathbf{x})_i}{\partial \mathbf{b}_j^o} = 1_{i=j}$

Backpropagation



Partial derivatives related to \mathbf{W}^o

- $\frac{\partial l}{\partial W_{i,j}^o} = \sum_k \frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_k} \frac{\partial \mathbf{z}^o(\mathbf{x})_k}{\partial W_{i,j}^o}$
- $\nabla_{\mathbf{W}^o} l = (\mathbf{f}(\mathbf{x}) - \mathbf{e}(y)) \cdot \mathbf{h}(\mathbf{x})^\top$

Backprop gradients

Compute activation gradients

- $\nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

Compute layer params gradients

- $\nabla_{\mathbf{W}^o} \mathbf{l} = \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l} \cdot \mathbf{h}(\mathbf{x})^\top$
- $\nabla_{\mathbf{b}^o} \mathbf{l} = \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l}$

Compute prev layer activation gradients

- $\nabla_{\mathbf{h}(\mathbf{x})} \mathbf{l} = \mathbf{W}^{o\top} \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l}$
- $\nabla_{\mathbf{z}^h(\mathbf{x})} \mathbf{l} = \nabla_{\mathbf{h}(\mathbf{x})} \mathbf{l} \odot \sigma'(\mathbf{z}^h(\mathbf{x}))$

Initialization and normalization

- Input data should be normalized to have approx. same range:
 - standardization or quantile normalization
- Initializing W^h and W^o :
 - Zero is a saddle point: no gradient, no learning
 - Constant init: hidden units collapse by symmetry
 - Solution: random init, ex: $w \sim \mathcal{N}(0, 0.01)$
 - Better inits: Xavier Glorot and Kaming He & orthogonal
- Biases can (should) be initialized to zero

SGD Learning Rate

- Very sensitive:
 - Too high → early plateau or even divergence
 - Too low → slow convergence
 - Try a large value first: $\eta = 0.1$ or even $\eta = 1$
 - Divide by 10 and retry in case of divergence
- Large constant LR prevents final convergence
 - multiply η_t by $\beta < 1$ after each update
 - or monitor validation loss and divide η_t by 2 or 10 when no progress
 - See [ReduceLROnPlateau](#) in Keras

Alternative optimizers

- SGD (with Nesterov momentum)
 - Simple to implement
 - Very sensitive to initial value of η
 - Need learning rate scheduling
- Adam: adaptive learning rate scale for each param
 - Global η set to 3e-4 often works well enough
 - Good default choice of optimizer (often)
- But well-tuned SGD with LR scheduling can generalize better than Adam (with naive l2 reg)...
- Promising stochastic second order methods: [K-FAC](#) and [Shampoo](#) can be used to accelerate training of very large models.

Next Class:

Support Vector Machines