

**Universität Stuttgart**  
KI – Institute for Artificial Intelligence  
Analytic Computing

# Machine Learning

## 11 Neural Networks Part 1



Prof. Dr. Steffen Staab

Nadeen Fatallah

Daniel Frank

Akram Sadat Hosseini

Jiaxin Pan

Osama Mohamed

Arvindh Arunbabu

Tim Schneider

Yi Wang

<https://www.ki.uni-stuttgart.de/>

# Learning Objectives

- Biological vs artificial neural networks
- Feed-forward networks
- XOR-Problem
- Activation functions
- Loss functions for binomial, multi-nomial distributions
- Computation graphs
- Backpropagation

# **1 Biological vs artificial neural networks**

# History of Machine Intelligence (Artificial Neural Networks)

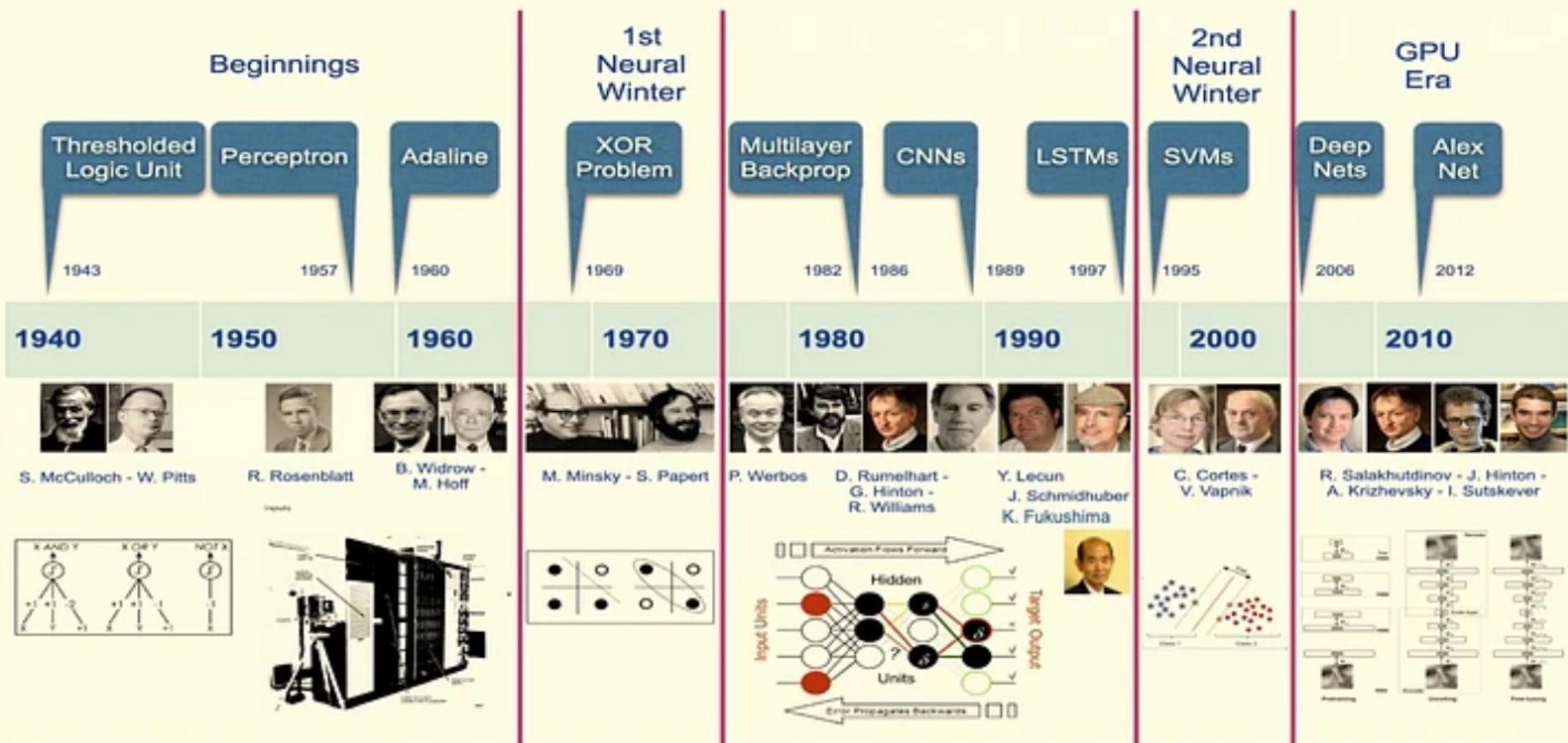
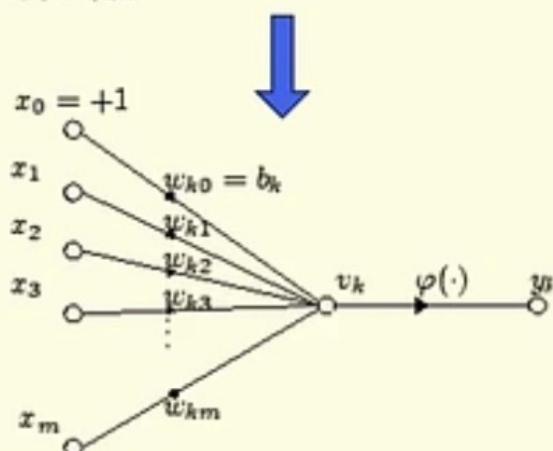
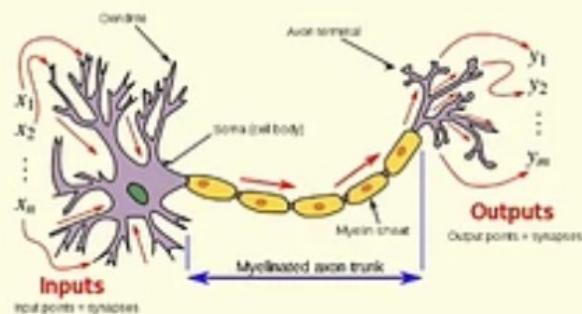


Figure courtesy of Professor Rene Vidal

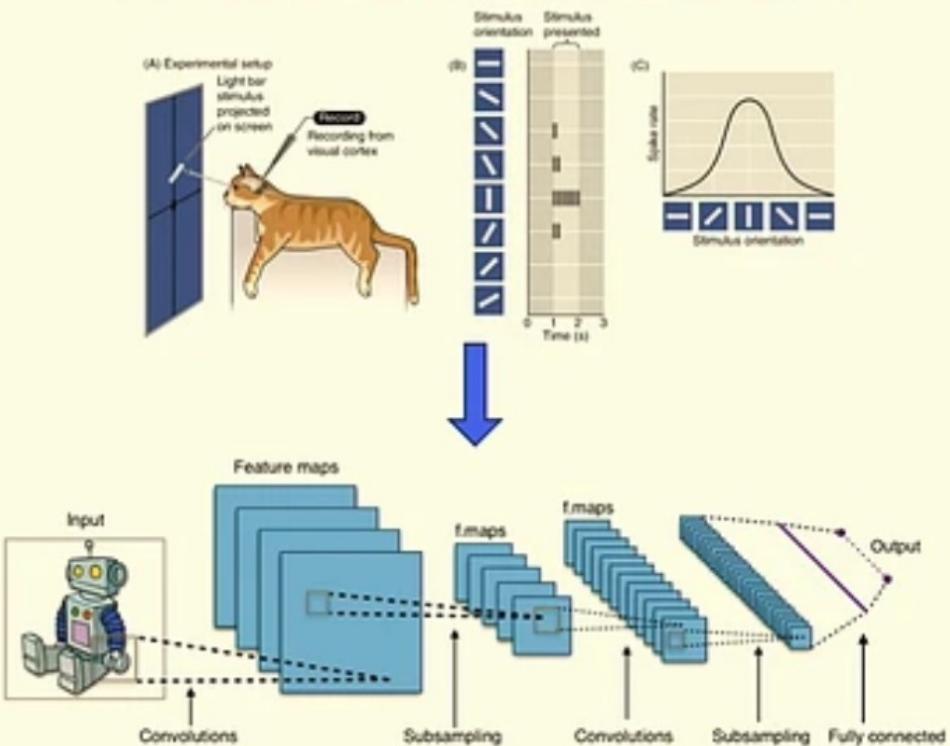
# Artificial Neurons and Neural Networks: Learn from Nature

Golgi and Cajal 1888 (1901 Nobel Prize)



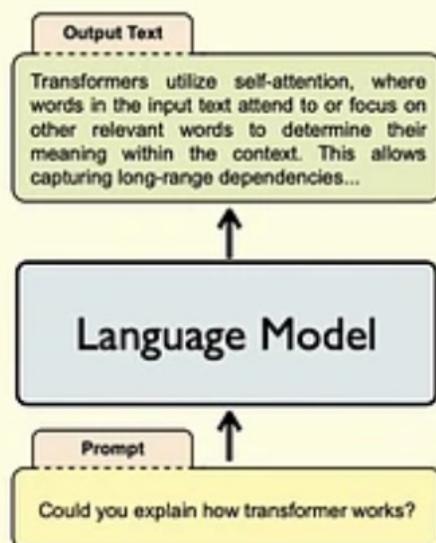
Warren McCulloch & Walter Pitts 1948

Hubel and Wiesel 1959 (1981 Nobel Prize)

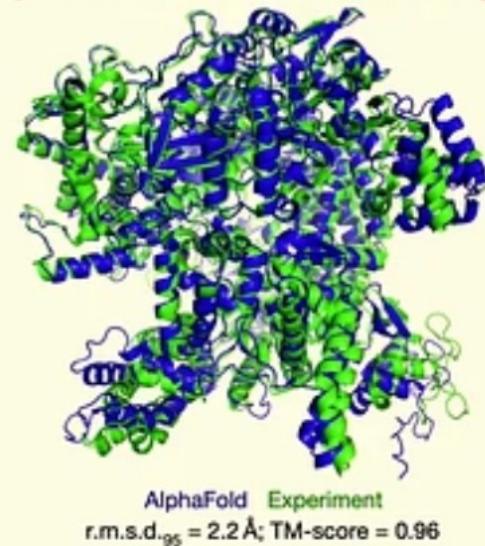


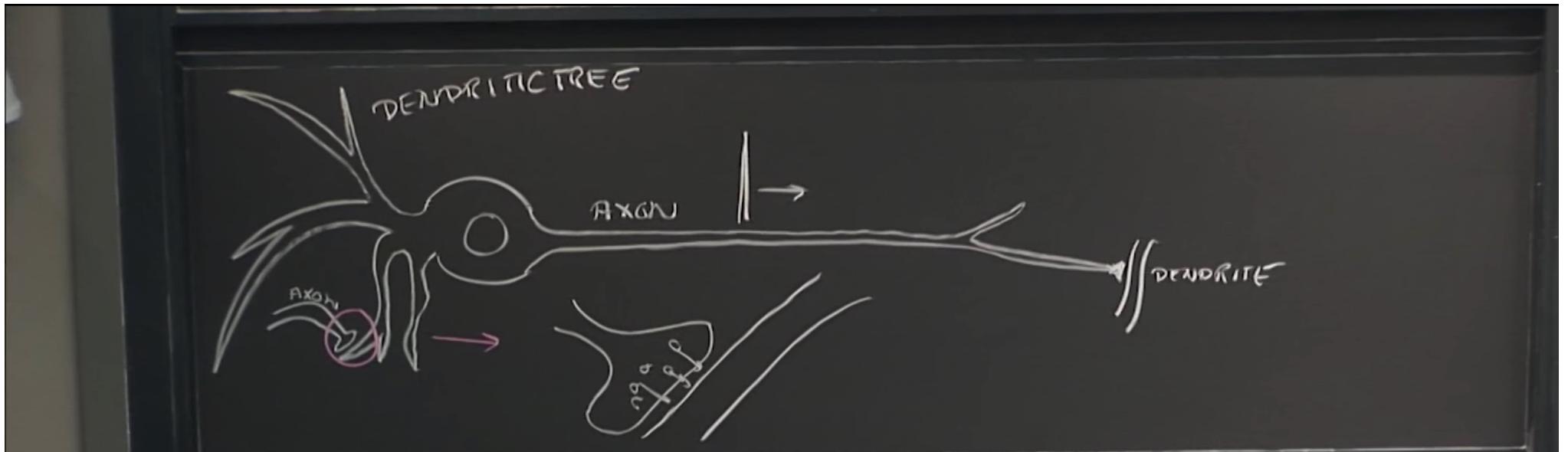
Fukushima 1980 & LeCun 1989 (Turing Award)

# Modern Evolution of Deep Neural Networks



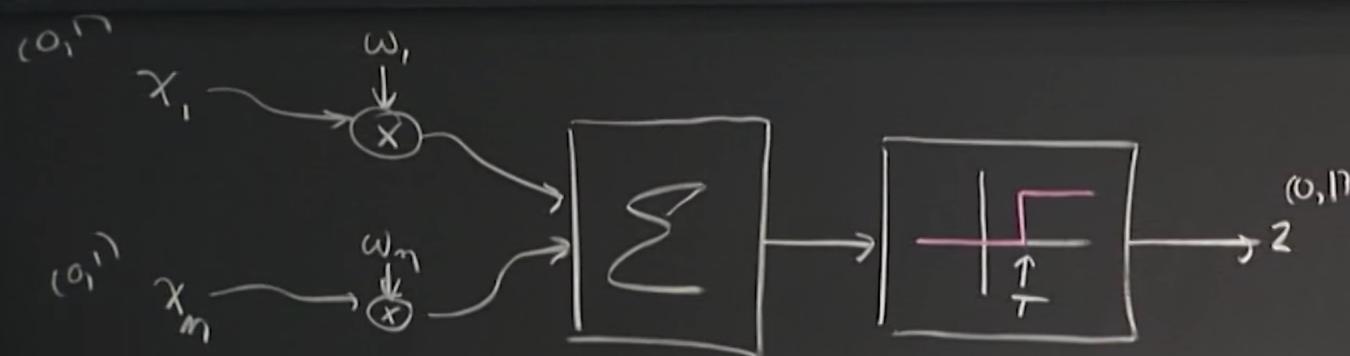
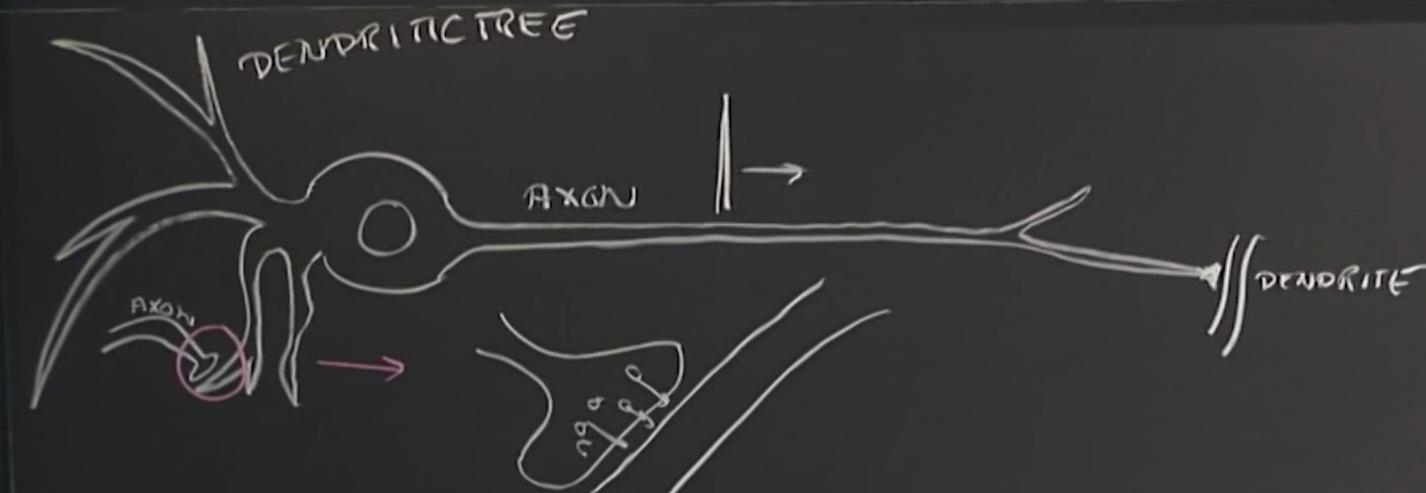
(2024 Nobel Prize Chemistry)





The working of natural neural networks ↑ motivated models for artificial neural networks.

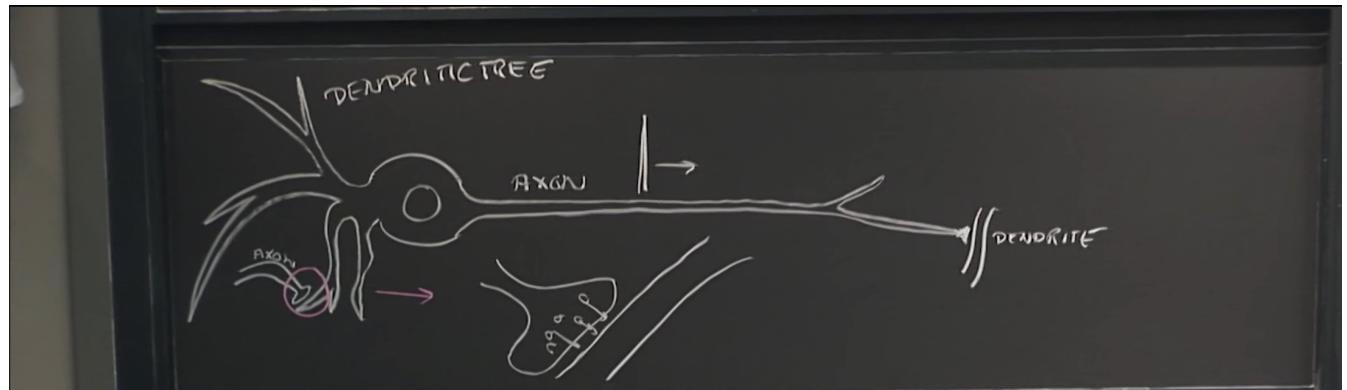
<https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-034-artificial-intelligence-fall-2010/lecture-videos/lecture-12a-neural-nets/>



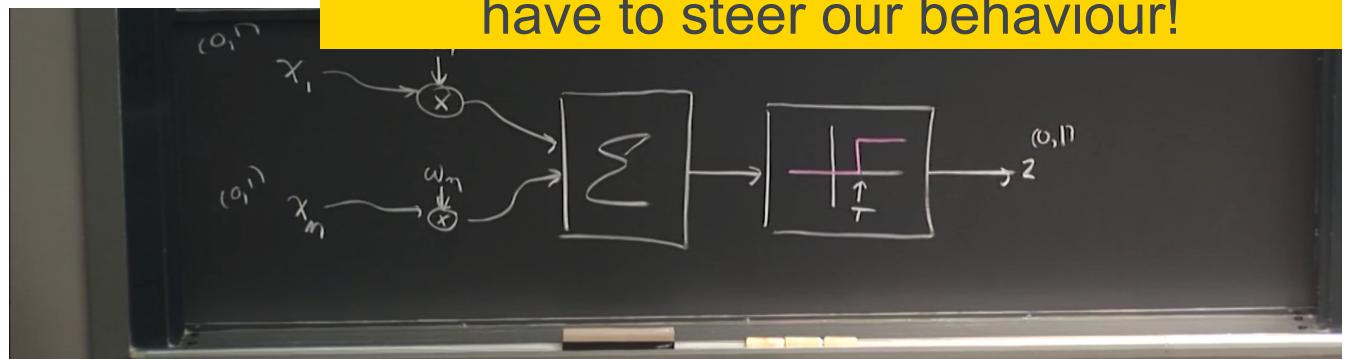
<https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-034-artificial-intelligence-fall-2010/lecture-videos/lecture-12a-neural-nets/>

# Differences between biological and (typical) artificial neural networks

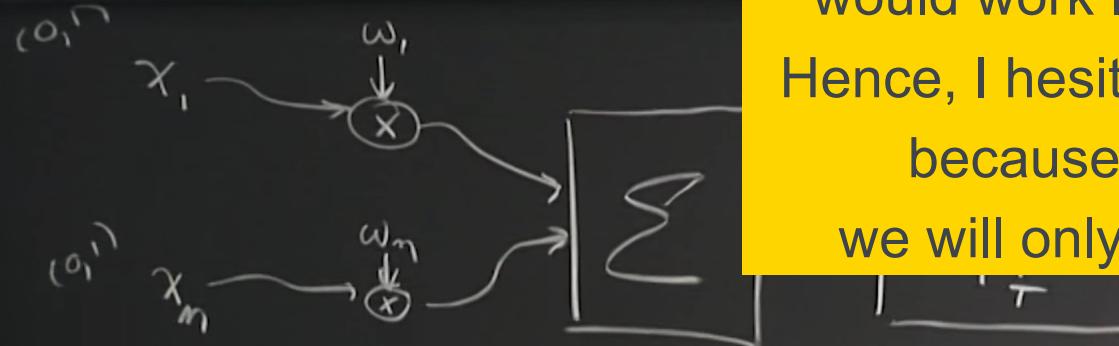
- Biological neuron
  - potential encoding
    - 100ms scale
  - temporal encoding
    - temporal coding
      - 1ms scale
    - differential equations
  - ~ Spiking neural network
- Artificial neural network (as typically used in ML)
  - non-linear statistical models



We could not walk if this encoding would have to steer our behaviour!

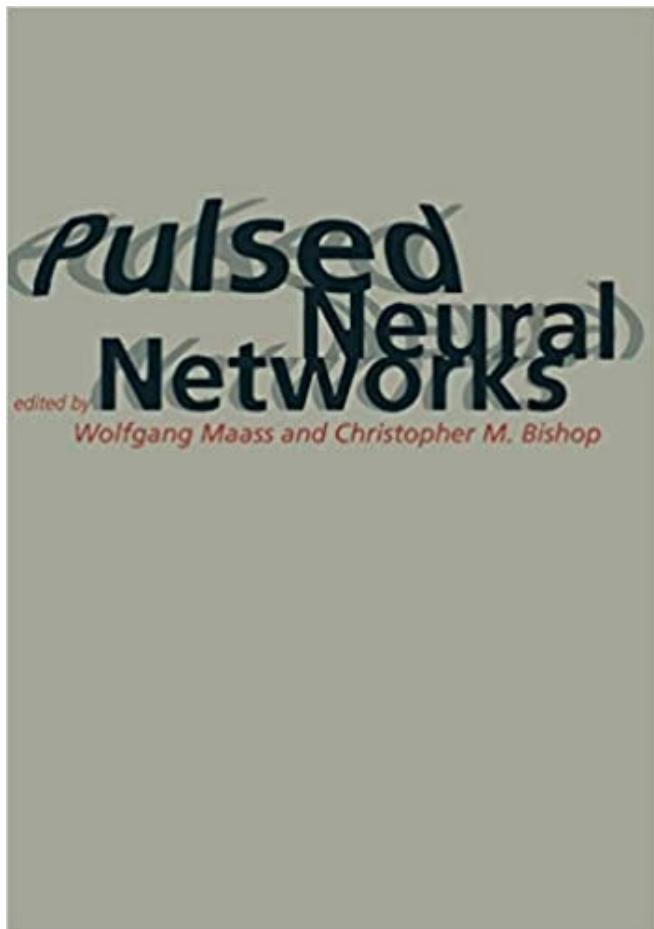


[https://en.wikipedia.org/wiki/Neural\\_coding](https://en.wikipedia.org/wiki/Neural_coding)



We could not walk if our neurons ↑  
would work like artificial neurons ↓ .  
Hence, I hesitate to motivate the latter  
because of the former and -  
we will only learn about the latter.

# Here, a book from 1999



From the book cover:

- Most practical applications of artificial neural networks are based on a computational model involving the propagation of continuous variables from one processing unit to the next. In recent years, data from neurobiological experiments have made it increasingly clear that biological neural networks, which communicate through pulses, **use the timing of the pulses to transmit information and perform computation**. This realization has stimulated significant research on pulsed neural networks, including theoretical analyses and model development, neurobiological modeling, and hardware implementation.
- This book presents the complete spectrum of current research in pulsed neural networks and includes the most important work from many of the key scientists in the field. Terrence J. Sejnowski's foreword, "Neural Pulse Coding," presents an overview of the topic. The first half of the book consists of longer tutorial articles spanning neurobiology, theory, algorithms, and hardware. The second half contains a larger number of shorter research chapters that present more advanced concepts. The contributors use consistent notation and terminology throughout the book.

**Unfortunately, pulsed neural networks have not been a success in technical applications until today. They are also not much researched (in comparison).**

# Neuromorphic computing

- IBM True North Chip

[https://en.wikipedia.org/wiki/Cognitive\\_computer#IBM\\_True\\_North\\_chip](https://en.wikipedia.org/wiki/Cognitive_computer#IBM_True_North_chip)

- ...

- TU Dresden SPINNAKER2 Chip  
(SPIking Neural Network Architecture)

Low power consumption!

- Liquid neural networks



## Liquid neural networks

- Simulating C. Elegans
- Using differential equations
- [https://www.quantamagazine.org/researchers-discover-a-more-flexible-approach-to-machine-learning-20230207/?fbclid=IwAR3qnPFC-EdF9jzMNu4XubABaM\\_S9HTKaia8ZjtTPM-xuRqyQLAyNIM2r8o](https://www.quantamagazine.org/researchers-discover-a-more-flexible-approach-to-machine-learning-20230207/?fbclid=IwAR3qnPFC-EdF9jzMNu4XubABaM_S9HTKaia8ZjtTPM-xuRqyQLAyNIM2r8o)

## **2 Reminder: Logistic regression**

# Logistic regression

## Model

- Composed of
  - Linear function:  $z = x^T \beta$
  - Non-linear **activation function**:  $y = \frac{e^z}{1+e^z}$
- Together:  $\hat{f}(x) = \frac{e^{x^T \beta}}{1+e^{x^T \beta}}$

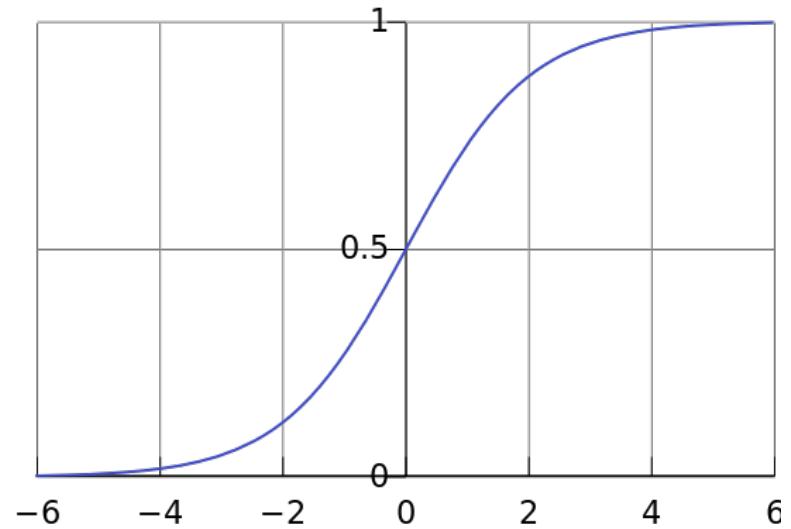
## Loss function

- Cross entropy and **regularization**

$$\text{loss}^{\text{logistic}}(\beta) = \sum_{i=1}^N H(\bar{y}_i, P(\cdot | X = x_i)) + \lambda \|\beta\|^2$$

## Solving

- adapting the weights, e.g. by gradient descent

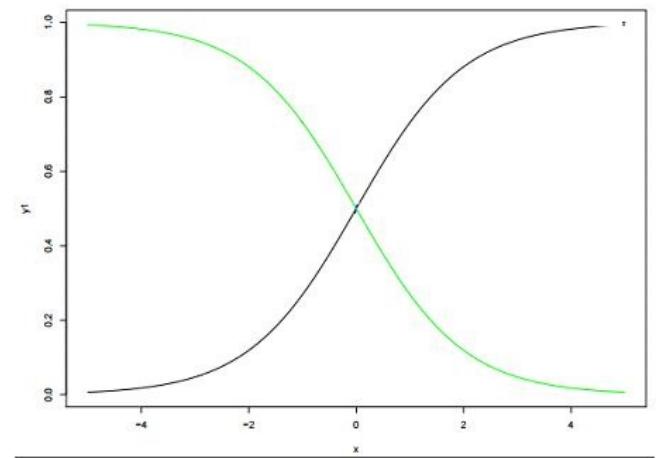


# Logistic regression for binary case

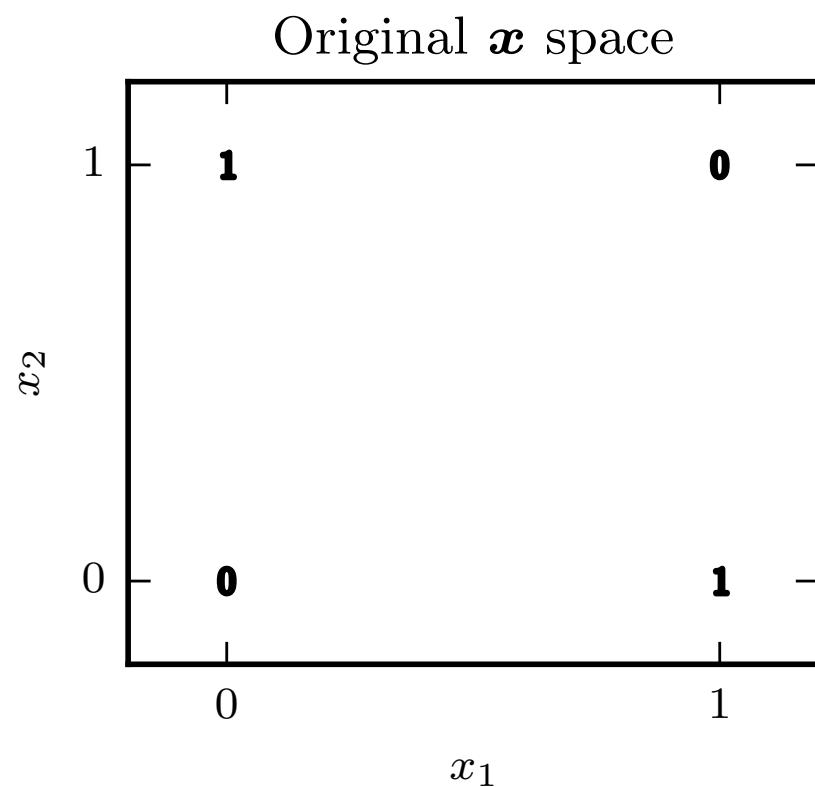
- Predict conditional probabilities  $\in [0,1]$

- 2-class example

$$\begin{aligned} \bullet P(G = 1|X = x) &= \frac{e^{x^T \beta}}{1+e^{x^T \beta}} \\ \bullet P(G = 2|X = x) &= \frac{1}{1+e^{x^T \beta}} \end{aligned} \quad \left. \vphantom{\frac{e^{x^T \beta}}{1+e^{x^T \beta}}} \right\} \text{Sum is 1}$$



## Problem: Learning XOR with logistic regression



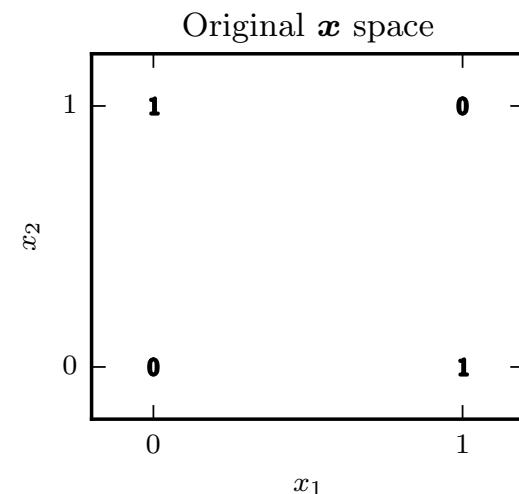
(Goodfellow 2017)

## Decision boundary is a hyperplane: Cannot learn XOR

$$\log \frac{P(G = 1|X = x)}{P(G = 2|X = x)} = \log \frac{\frac{e^{x^T \beta}}{1 + e^{x^T \beta}}}{\frac{1}{1 + e^{x^T \beta}}} = \log e^{x^T \beta} = x^T \beta$$

Dealing with non-linear separability

1. Kernelization - also works for logistic regression!
2. Composition of non-linear functions



# **3 Feed-Forward Networks and Activation Functions**

## Networks: Composing functions

Three layer network:

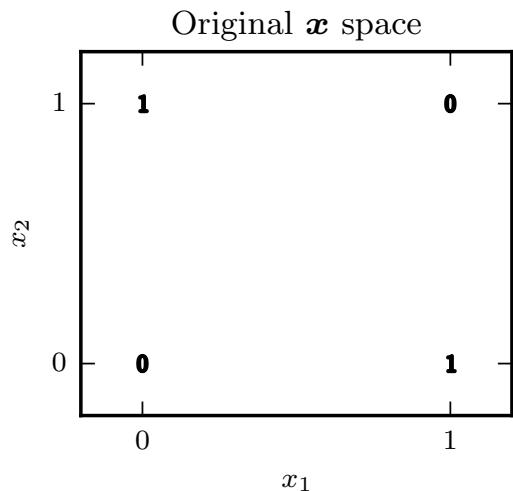
$$f(x) = f^{(3)} \left( f^{(2)} \left( f^{(1)}(x) \right) \right)$$

First layer:  $f^{(1)}$

Second layer:  $f^{(2)}$

Last layer, output layer:  $f^{(3)} = \hat{y}$

# What about composition of linear functions?



Find solution

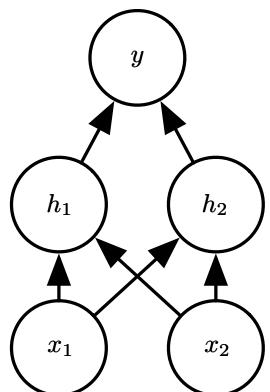
$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = f^2(f^1(\mathbf{x}))^T$$

With

$$\begin{aligned} f^1(\mathbf{x}; \mathbf{W}, \mathbf{c}) &= \mathbf{h} \\ f^2(\mathbf{h}; \mathbf{w}, b) &= y \end{aligned}$$

Then

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{W}^T \mathbf{x}$$



A composition of linear functions gives a linear function

→ cannot solve XOR learning problem

# One Layer in a Feed-Forward Network

A linear combination of non-linear basis functions (features):

$$f^{(i)}(x) = w_{i,0} + \sum_{j=1}^{l_i} w_{i,j} \cdot \phi(x; \theta_i) = \phi(x; \theta)^T w_i$$

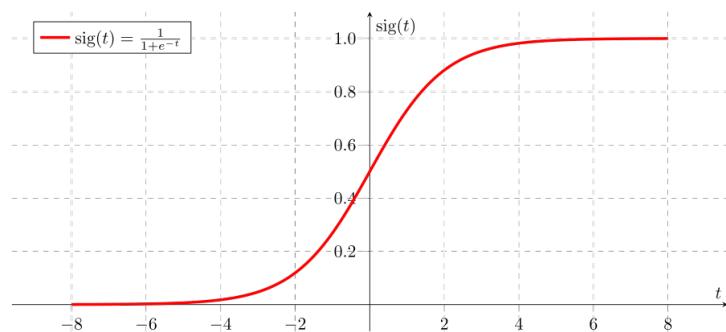
with  $l_i$  being **width** of layer  $i$ .

How to pick  $\phi(x; \theta_i)$ ?

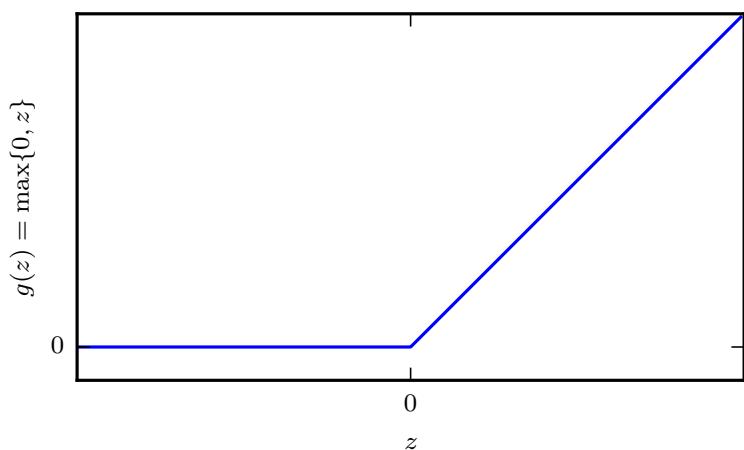
- Core idea:  $\phi(x; \theta_i) = g(x^T \theta_i)$ 
  - learnable linear component:  $x^T \theta_i$
  - fixed non-linear activation function:  $g(z)$

# Activation Functions

Sigmoid (cf logistic regression)



ReLU – Rectified Linear Unit



Sigmoid:

$$g(z) = \sigma(z) = \frac{e^z}{1+e^z}$$

ReLU:

$$g(z) = \max\{0, z\}$$

Leaky ReLU:  
(never saturates fully)

$$g_\alpha(z) = \max\{\alpha z, z\},  
e.g. \alpha = 0.1$$

Maxout unit:

$$g(z) = \max_j z_j$$

(largest value of group)

Absolute value rectification:  $g(z) = |z|$   
(e.g. for object recognition)

Hyperbolic tangent:

$$g(z) = \tanh(z)$$

Radial Basis Function:

$$g(z) = e^{-\frac{1}{2}\|z-c_j\|^2}$$

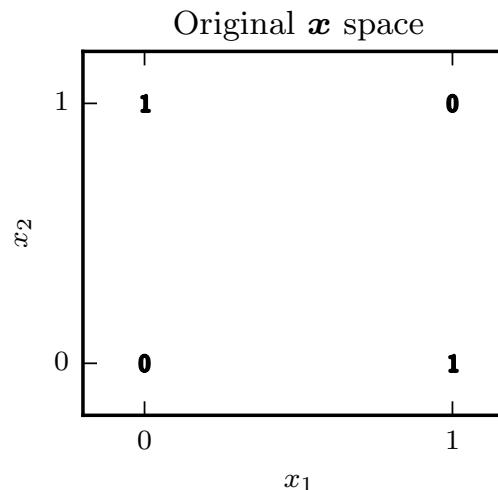
- + Simple
- + Efficient to compute with
- + Non-linear
- Not differentiable at 0

## Possible Solution to the XOR problem

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

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

$$Y^T = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$



## Possible Solution to the XOR problem

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

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

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

$$w^T X = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \end{bmatrix}$$

$$c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

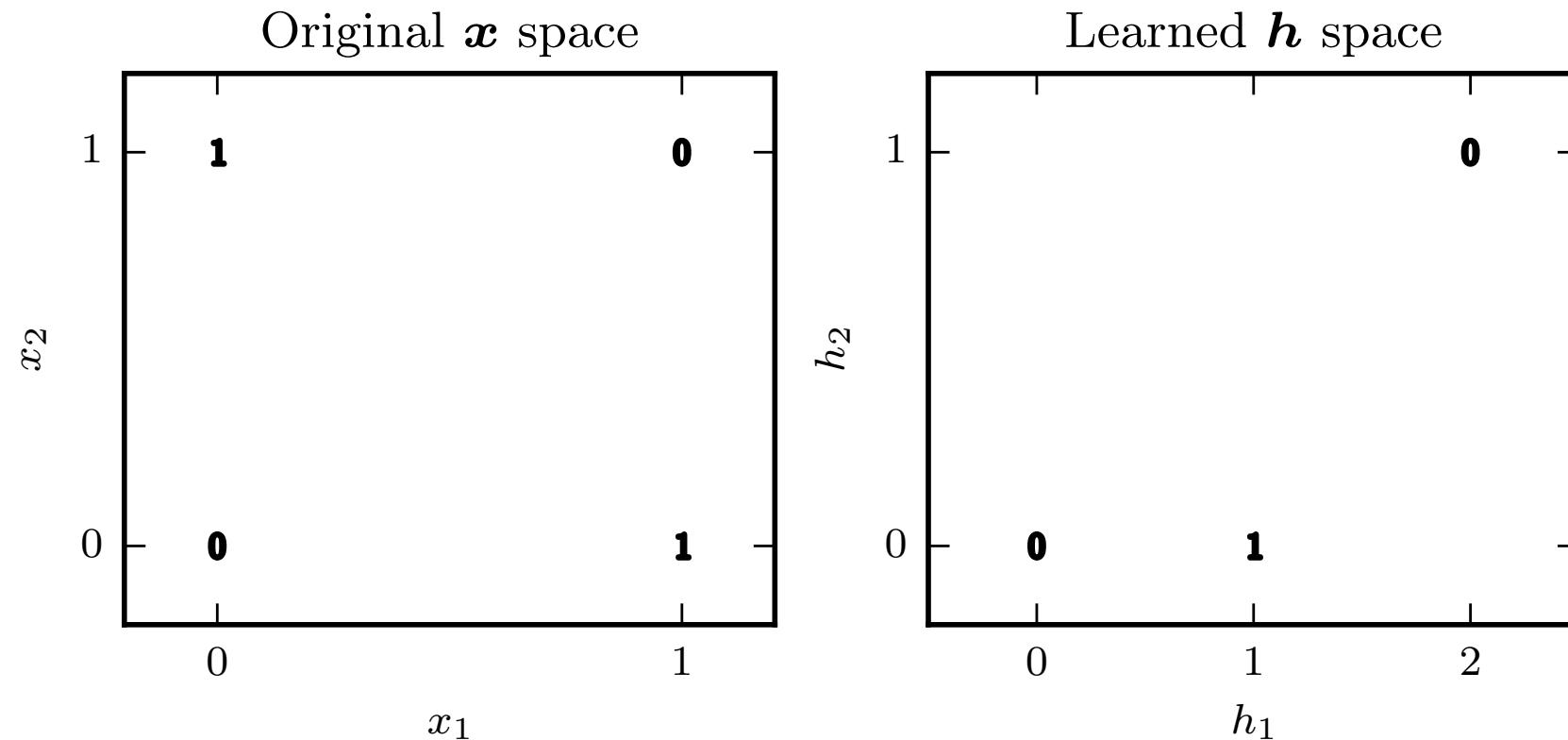
$$\begin{aligned} w^T X + [c & c & c & c] \\ = \begin{bmatrix} 0 & 1 & 1 & 2 \\ -1 & 0 & 0 & 1 \end{bmatrix} \end{aligned}$$

ReLU (pointwise)  
 $\begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix}$

$$w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$\begin{aligned} [1 & -2] \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix} \\ = [0 & 1 & 1 & 0] \end{aligned}$$

## Solving XOR problem



# Universal approximation theorem

**Theorem:** Let  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  be a nonconstant, bounded, and monotonically-increasing continuous function.

Let  $C$  denote the space of continuous function on  $[0,1]^m$ .

Given **any**  $\varepsilon > 0$  and **any** function  $f \in C$ ,

there exist  $N \in \mathbb{N}$ ,  $v_i, b_i \in \mathbb{R}$ ,  $w_i \in \mathbb{R}^m$ , with  $i = 1, \dots, N$ ,  
such that

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i)$$

and

$$|F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon$$

for all  $\mathbf{x} \in [0,1]^m$ . □

On compact spaces, simple, sufficiently broad neural networks with one type of non-linear activation function approximate any continuous function arbitrarily well.

BUT: This is **not** a practical proposition!  $N$  might have to be exponential in number of data points.

# Summary: One Layer Feed-forward Network

## Model

- Composed of

- Linear function:  $z = x^T W + c$

- Non-linear **activation function**  $g(z)$ :

$$g(z) = \text{ReLU}(z), g(z) = \sigma(z) = \frac{e^z}{1 + e^z}, g(z) = \tanh(z), g(z) = e^{-\frac{1}{2}\|z - c_j\|^2}, \dots$$

- Together:  $\hat{f}(x) = w^T g(x^T W + c) + b$

- If we capture the „bias unit“ and write  $x^T = (1 \ x_1 \ \dots \ x_d)$ , we can (again) simplify this to  $\hat{f}(x) = w^T g(x^T W)$

## Loss function ??

## Solving ???

## Chain-based architecture (fully connected)

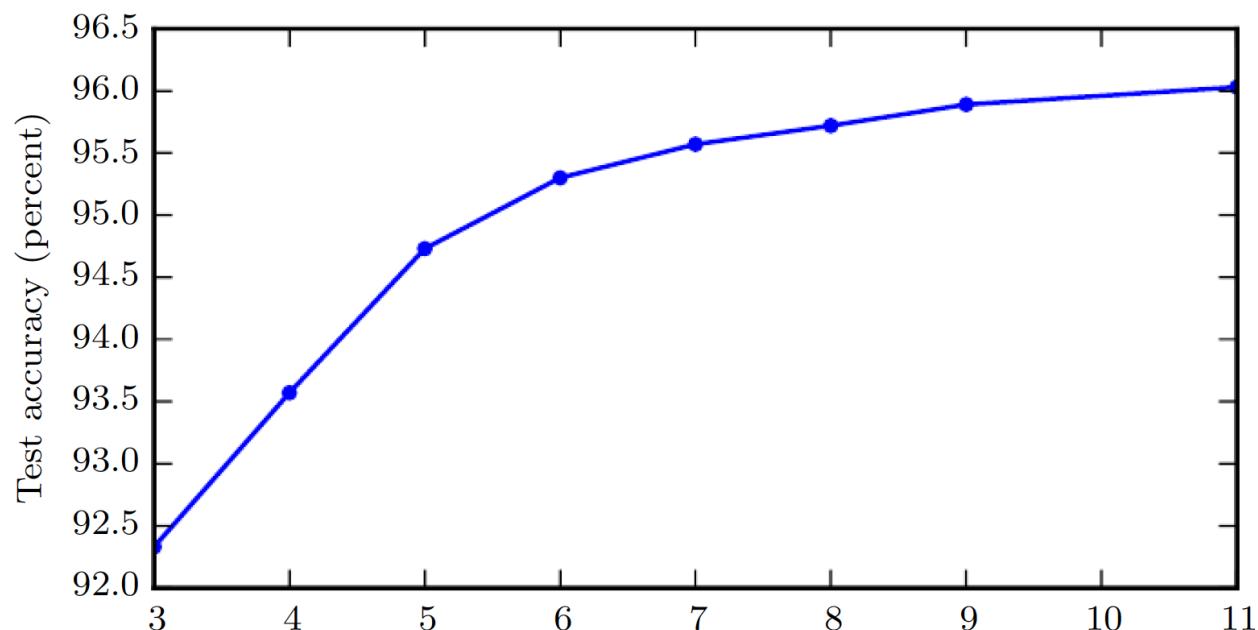
- Input layer (not counted):  $x \in \mathbb{R}^{d+1}, x = h^{(0)}$
- First layer:  $h^{(1)} = g^{(1)}(W^{(1)T}x + b^{(1)})$
- Second layer:  $h^{(2)} = g^{(2)}(W^{(2)T}h^{(1)} + b^{(2)})$
- ...
- n-th layer:  $h^{(n)} = g^{(n)}(W^{(n)T}h^{(n-1)} + b^{(n)}) = \hat{y}$

## Design considerations

- Depth of network
- Width of network  $l_i$  at each layer  $i$
- $W^i \in \mathbb{R}^{l_i \times \mathbb{R}^{l_{i-1}}}$

## Effect of depth

Experiment on recognizing multidigit numbers with increasing depth (convolutional\* with 1 or 2 fully connected layers)

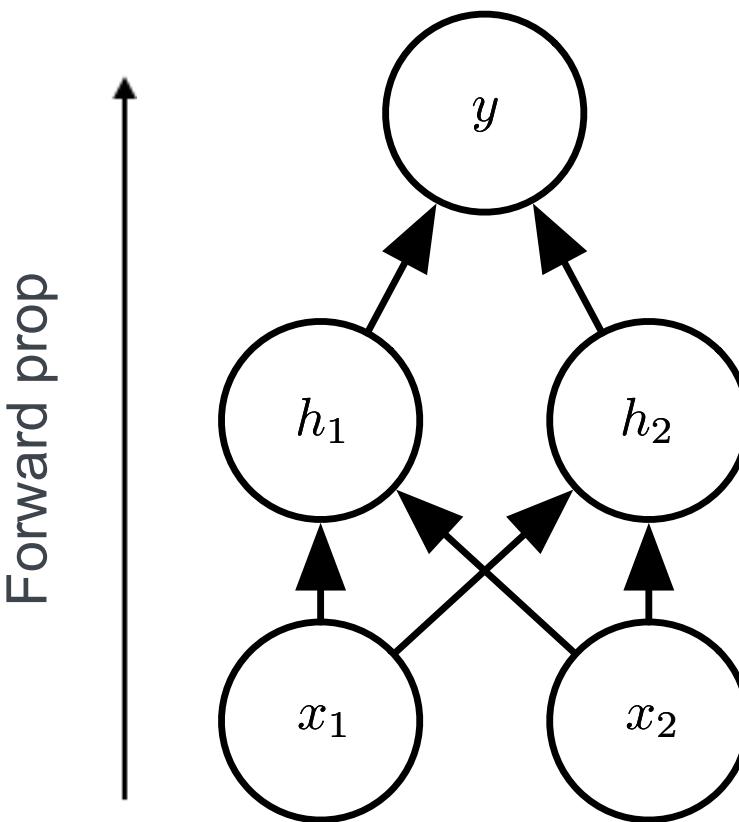


\*Convolutional network to be explained in next lecture

# 4 Backpropagation

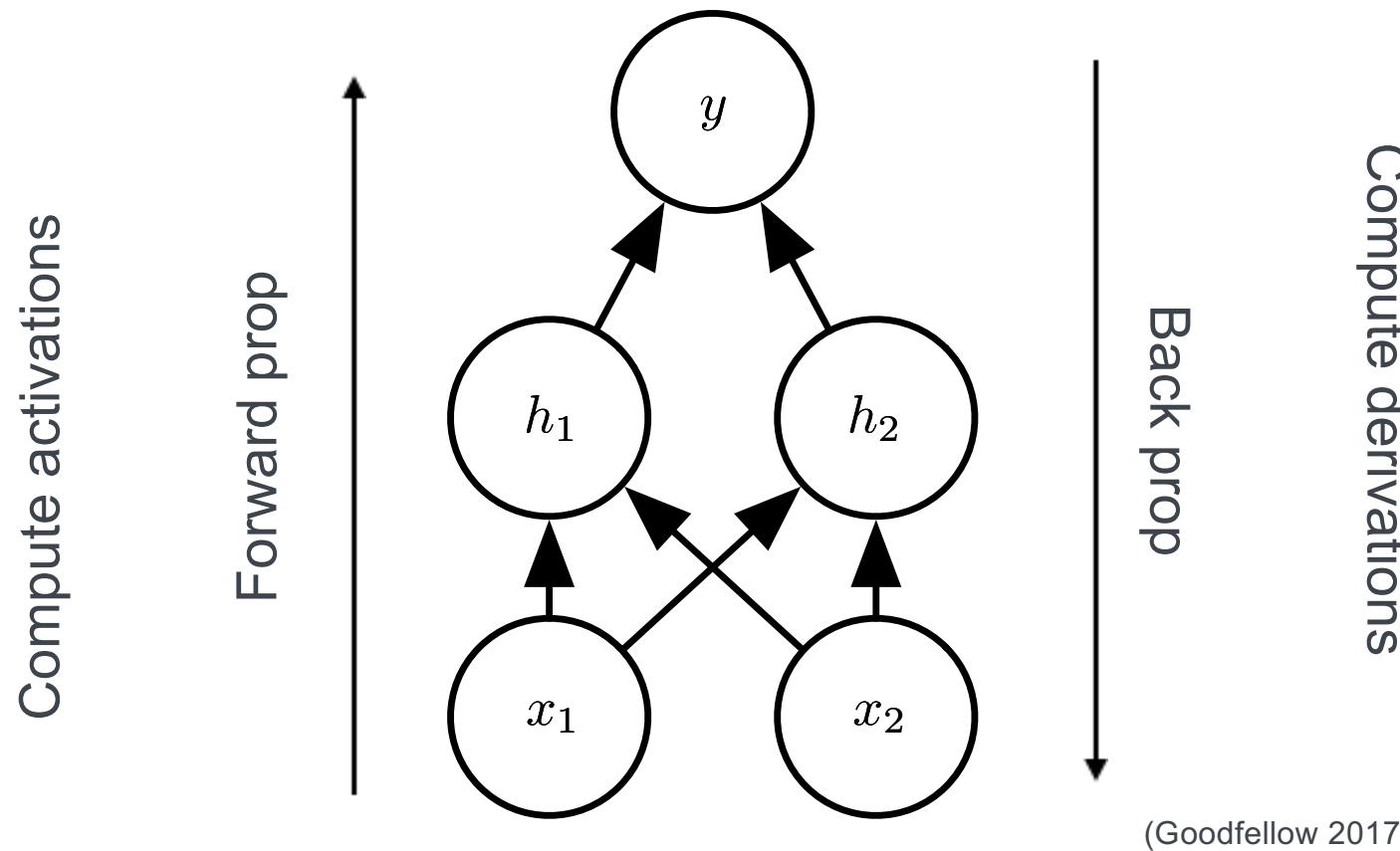
# Forward Propagation

Compute activations



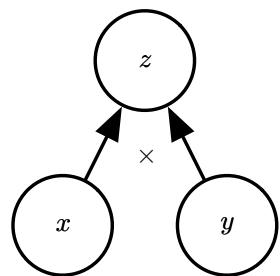
# Back-Propagation

Compute loss  $L(\theta)$  considering  $y$  (data) vs  $\hat{y}$  (output activation/prediction)



(Goodfellow 2017)

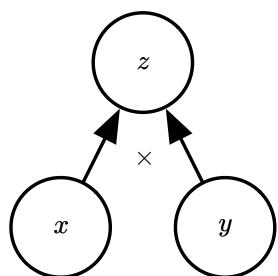
# Language for representing the backprop problem: Computation Graphs



(a)

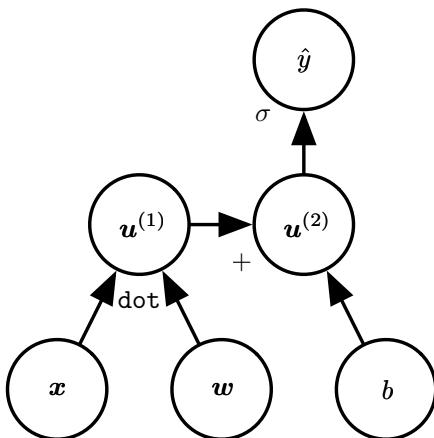
Multiplication

# Language for representing the backprop problem: Computation Graphs



(a)

Multiplication



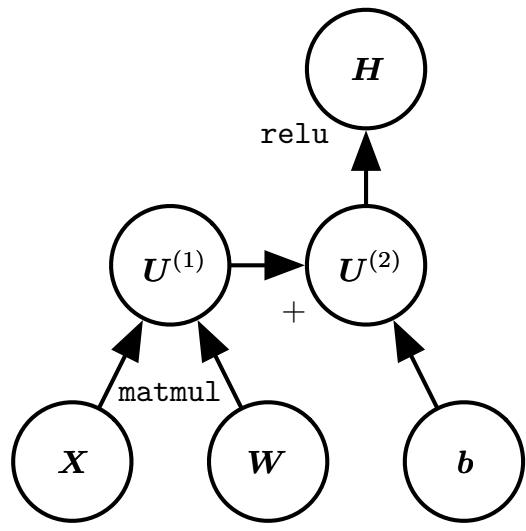
(b)

Logistic regression

Syntax elements:

- **Nodes represent variables**
- (multiple) edges with labels represent **operations**

# Language for representing the backprop problem: Computation Graphs



(c)

ReLU Layer

Syntax elements:

- **Nodes** represent **variables**
- (multiple) edges with labels represent **operations**

# Some notational conventions for computation graph

$n$  many values of nodes in the network

- topologically sorted

Example:

$$u_1 \dots u_6$$

Input layer with values for  $k$  nodes:

$$u_1 \dots u_k$$

Set of values of parents of node  $i$ :

$$\mathbb{A}^{(i)}$$

Operator for node  $i$ :

$$f^{(i)}$$

Output layer with  $l$  nodes:

$$u_{n-l+1} \dots u_n$$

$$u_1 = x_1, u_2 = w, u_3 = b$$

$$u_4 = u^{(1)}, u_5 = u^{(2)}$$

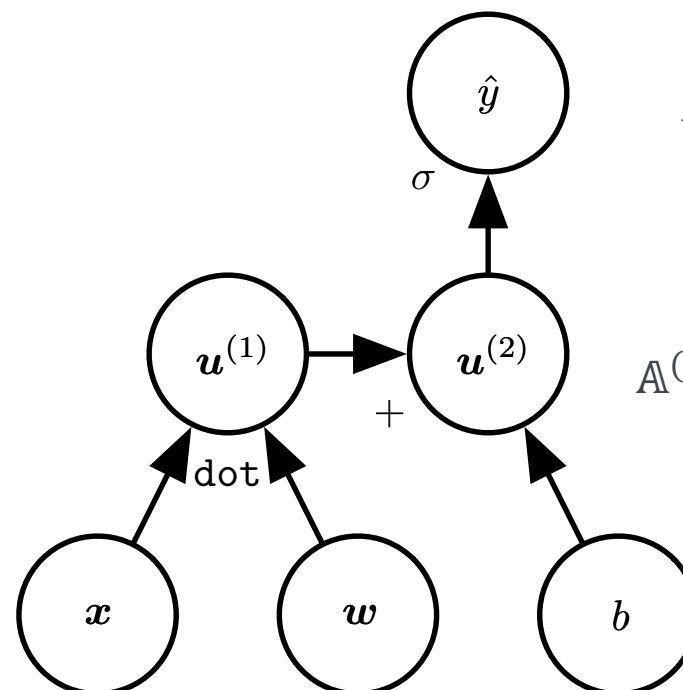
$$\mathbb{A}^{(4)} = \{u_1, u_2\}$$

$$\mathbb{A}^{(5)} = \{u_3, u_4\} = \{b, u^{(1)}\}$$

$$\mathbb{A}^{(6)} = \{u_5\}$$

$$f^{(4)} = \text{dot}, f^{(5)} = +$$

$$f^{(6)} = \sigma$$



$$u_{6_{37}} = \hat{y}$$

# Forward propagation in computation graph

$n$  many values of nodes in the network

Input layer with values for  $k$  nodes:

$$u_1 \dots u_k$$

Set of values of parents of node  $i$ :

$$\mathbb{A}^{(i)}$$

Operator for node  $i$ :

$$f^{(i)}$$

Output layer with 1 nodes:

$$u_{n-l+1} \dots u_n$$

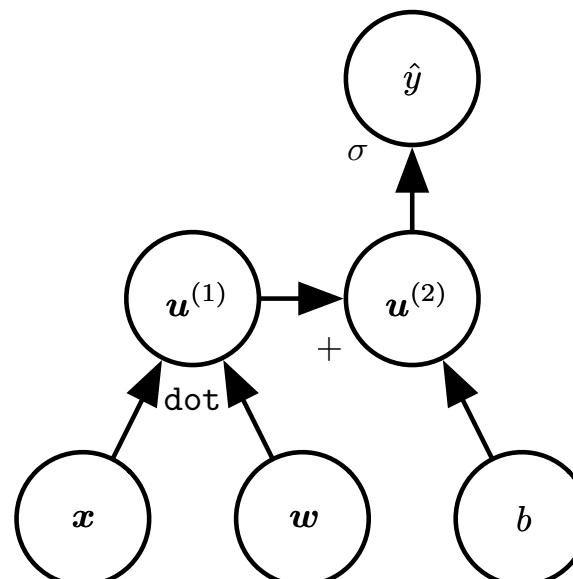
**for**  $i = k + 1 \dots n$  **do**

$$\mathbb{A}^{(i)} \leftarrow \{u_j \mid j \in \text{parent}(u_i)\};$$

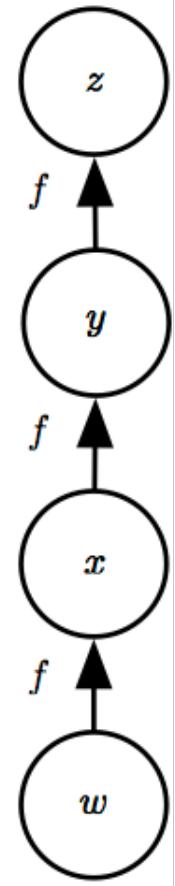
$$u_i \leftarrow f^{(i)}(\mathbb{A}^{(i)});$$

**end for**

**return**  $u_n$



## Some ideas of backpropagation in a simple example



$$\frac{\partial z}{\partial w} \tag{6.50}$$

$$= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \tag{6.51}$$

$$= f'(y) f'(x) f'(w) \tag{6.52}$$

$$= f'(f(f(w))) f'(f(w)) f'(w) \tag{6.53}$$

Back-prop avoids computing this twice

# Backpropagation in Tiny Example Computation Graph

Input

$$x = 1, y = 2; \text{target } z = 3$$

Forward result:

$$\hat{z} = 2$$

Increased  
 $\hat{z}$  reduces  
error

Assume MSE as loss function:

$$J = L(\hat{z}, z) = \frac{1}{2}(\hat{z} - z)^2 = \frac{1}{2}$$

Gradient on loss function  $L$ :

$$\begin{aligned}\nabla_{\hat{z}} J &= \nabla_{\hat{z}} L(\hat{z}, z) = 2 \frac{1}{2}(\hat{z} - z) \\ &= (\hat{z} - z) = \hat{z} - 3 = -1\end{aligned}$$

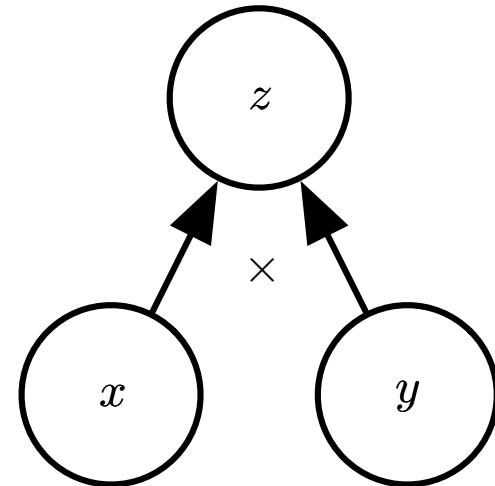
Backpropagation

$$\nabla_{(x,y)} J = \begin{pmatrix} -2 \\ -1 \end{pmatrix}$$

Partial derivatives

$$\nabla_x J = \nabla_{\hat{z}} J \cdot \nabla_x \hat{z} = -1 \cdot \frac{\partial(x \cdot y)}{\partial x} = -y = -2$$

$$\nabla_y J = \nabla_{\hat{z}} J \cdot \nabla_y \hat{z} = -1 \cdot \frac{\partial(x \cdot y)}{\partial y} = -x = -1$$



Conclusion:  
If  $y$  was a weight it should be increased  
to bring the error down

$$y^{new} = y - \eta(-1)$$

# Larger Example - Forward: Fully connected feed-forward network

**Require:** Network depth,  $l$

**Require:**  $W^{(i)}, i \in \{1, \dots, l\}$ , the weight matrices of the model

**Require:**  $b^{(i)}, i \in \{1, \dots, l\}$ , the bias parameters of the model

**Require:**  $x$ , the input to process

**Require:**  $y$ , the target output

$$h^{(0)} = x$$

**for**  $k = 1, \dots, l$  **do**

$$a^{(k)} = b^{(k)} + W^{(k)}h^{(k-1)}$$

$$h^{(k)} = f(a^{(k)})$$

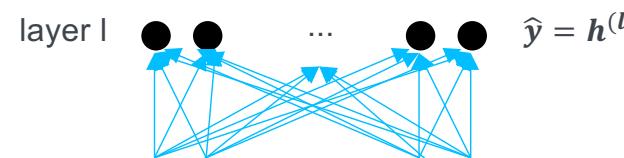
**end for**

$$\hat{y} = h^{(l)}$$

$$J = L(\hat{y}, y) + \lambda \Omega(\theta)$$

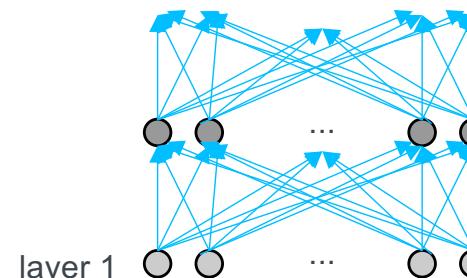
regularization

(Goodfellow 2017)



$$\hat{y} = h^{(l)}$$

.....



$$\begin{aligned} h^{(1)} &= f(a^{(1)}) = \\ &= f(b^{(1)} + W^{(1)}h^{(0)}) \end{aligned}$$

$$h^{(0)} = x$$

---

**Algorithm 6.3** Forward propagation through a typical deep neural network and the computation of the cost function. The loss  $L(\hat{\mathbf{y}}, \mathbf{y})$  depends on the output  $\hat{\mathbf{y}}$  and on the target  $\mathbf{y}$  (see section 6.2.1.1 for examples of loss functions). To obtain the total cost  $J$ , the loss may be added to a regularizer  $\Omega(\theta)$ , where  $\theta$  contains all the parameters (weights and biases). Algorithm 6.4 shows how to compute gradients of  $J$  with respect to parameters  $\mathbf{W}$  and  $\mathbf{b}$ . For simplicity, this demonstration uses only a single input example  $\mathbf{x}$ . Practical applications should use a minibatch. See section 6.5.7 for a more realistic demonstration.

---

**Require:** Network depth,  $l$

**Require:**  $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$ , the weight matrices of the model

**Require:**  $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$ , the bias parameters of the model

**Require:**  $\mathbf{x}$ , the input to process

**Require:**  $\mathbf{y}$ , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

**for**  $k = 1, \dots, l$  **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

**end for**

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

regularization

# Larger Example - Backward: Fully connected feed-forward network

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

**for**  $k = l, l-1, \dots, 1$  **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if  $f$  is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases  
where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

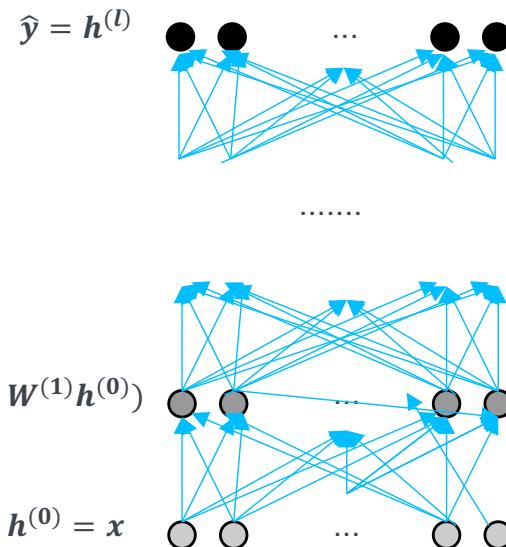
$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients  
w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

$$\mathbf{h}^{(1)} = f(\mathbf{a}^{(1)}) = f(\mathbf{b}^{(1)} + \mathbf{W}^{(1)} \mathbf{h}^{(0)})$$

$$\mathbf{h}^{(0)} = \mathbf{x}$$



---

**Algorithm 6.4** Backward computation for the deep neural network of algorithm 6.3, which uses, in addition to the input  $\mathbf{x}$ , a target  $\mathbf{y}$ . This computation yields the gradients on the activations  $\mathbf{a}^{(k)}$  for each layer  $k$ , starting from the output layer and going backwards to the first hidden layer. From these gradients, which can be interpreted as an indication of how each layer's output should change to reduce error, one can obtain the gradient on the parameters of each layer. The gradients on weights and biases can be immediately used as part of a stochastic gradient update (performing the update right after the gradients have been computed) or used with other gradient-based optimization methods.

---

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

**for**  $k = l, l - 1, \dots, 1$  **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if  $f$  is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

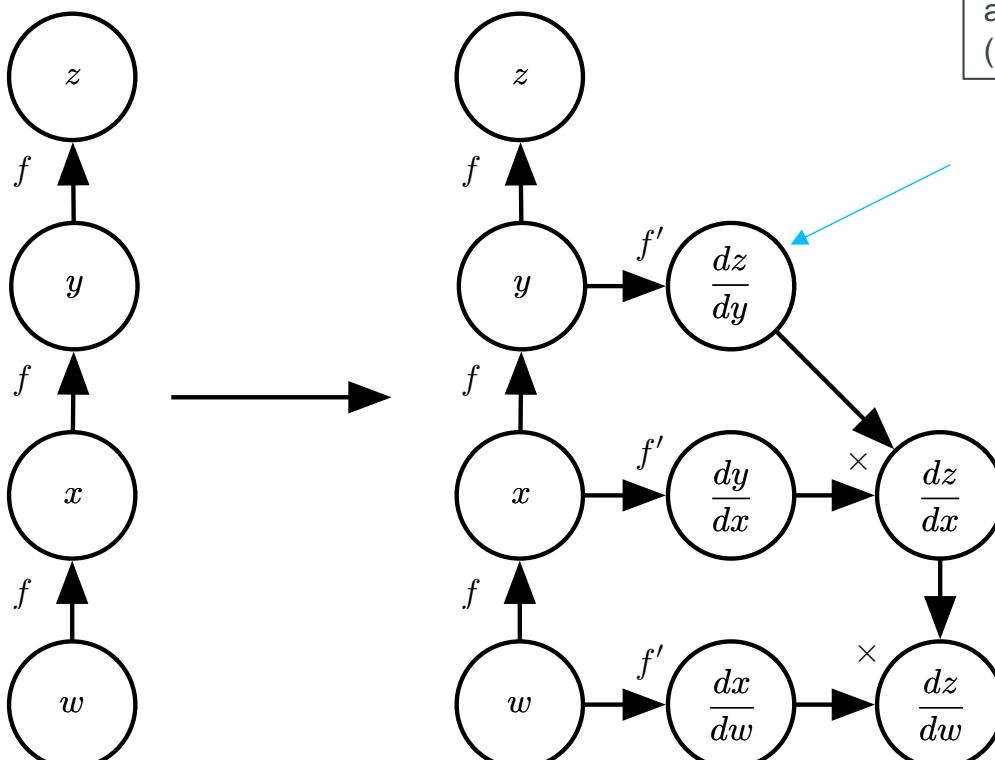
Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

**end for**

---

# Generalization 1: Symbol-to-Symbol derivatives



For each node representing a value, add a new node for (partial) derivatives

Some derivative nodes might be needed twice, but are computed only once, **but not in this example**

Backpropagation graph

## Generalization 2: Avoid double computation

**Algorithm 6.6** The inner loop subroutine `build_grad(V, G, G', grad_table)` of the back-propagation algorithm, called by the back-propagation algorithm defined in algorithm 6.5.

Require:  $V$ , the variable whose gradient should be added to  $G$  and `grad_table`

Require:  $G$ , the graph to modify

Require:  $G'$ , the restriction of  $G$  to nodes that participate in the gradient

Require: `grad_table`, a data structure mapping nodes to their gradients

```
if  $V$  is in grad_table then
    Return grad_table[V]
end if
 $i \leftarrow 1$ 
for  $C$  in get_consumers(V, G') do
     $op \leftarrow \text{get\_operation}(C)$ 
     $D \leftarrow \text{build\_grad}(C, G, G', \text{grad\_table})$ 
     $\mathbf{G}^{(i)} \leftarrow op.bprop(\text{get\_inputs}(C, G'), V, D)$ 
     $i \leftarrow i + 1$ 
end for
 $\mathbf{G} \leftarrow \sum_i \mathbf{G}^{(i)}$ 
grad_table[V] = G
Insert  $\mathbf{G}$  and the operations creating it into  $G$ 
Return  $\mathbf{G}$ 
```

} Dynamic programming      }

} If available, look up,  
else compute

We will not ask for  
details of this slide  
in the exam

## 5 Example application

# The optimization of a deep neural network

- is a rather ill-posed problem,
- has many degrees of freedom (compared to linear models),



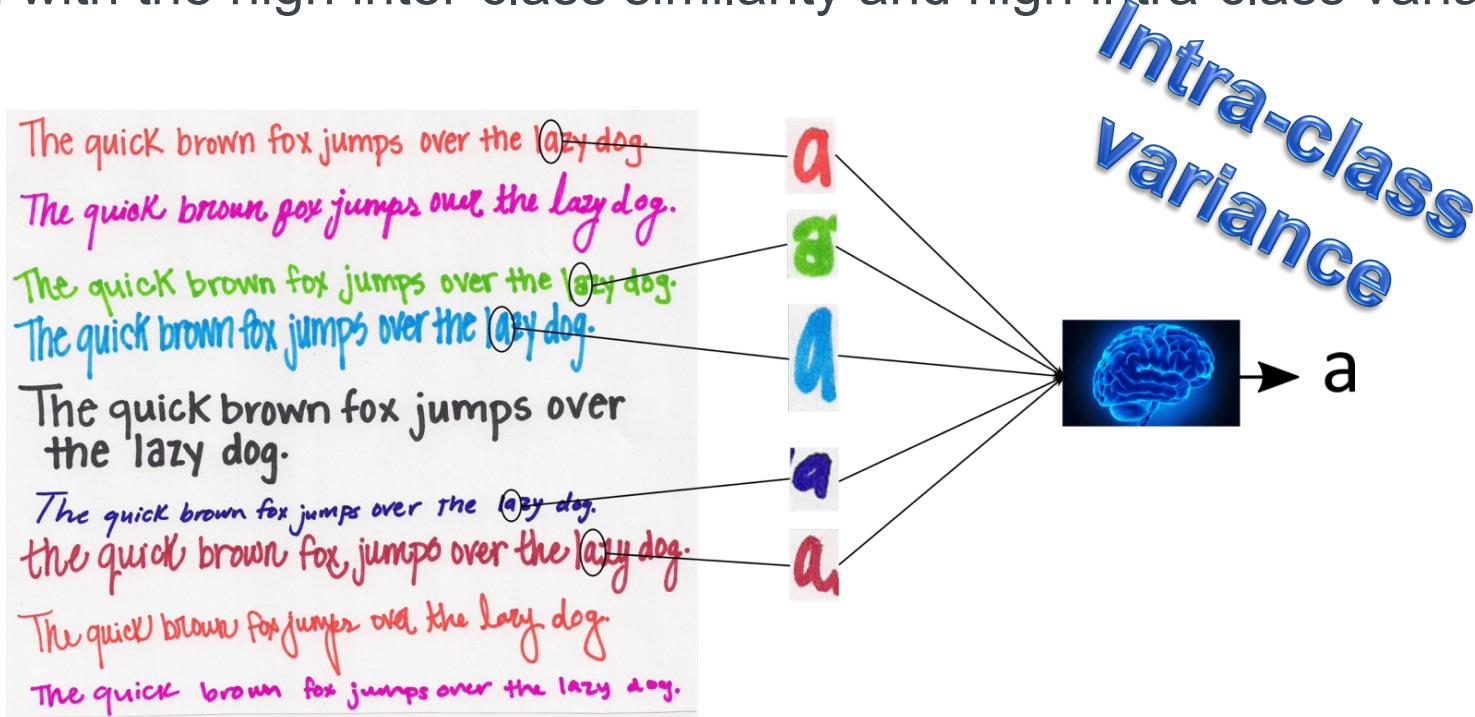
- Deep Neural Network especially useful if there is no previous knowledge about the data generating process

- image
- audio
- text

In these areas deep neural networks outperform other areas by far

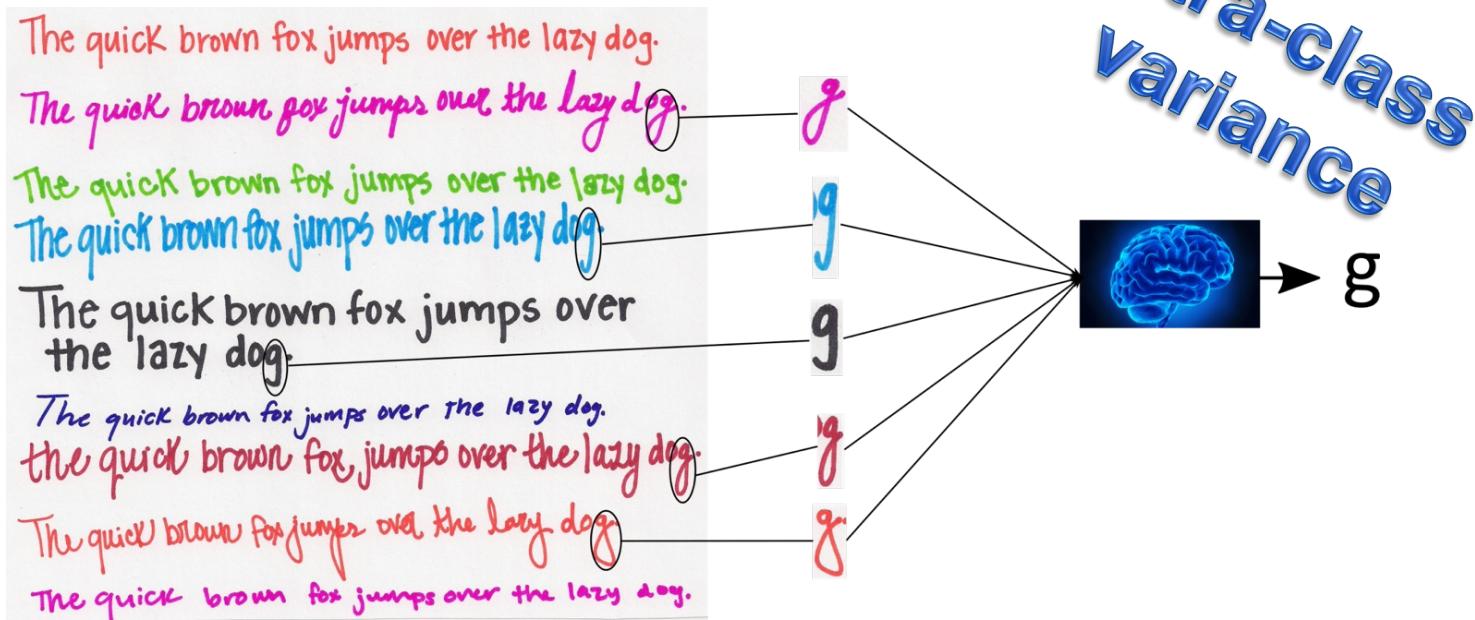
# Neural Networks

- Example: Handwritten characters classification
    - It is very simple for human beings to understand handwritten texts, even with the high inter-class similarity and high intra-class variance.



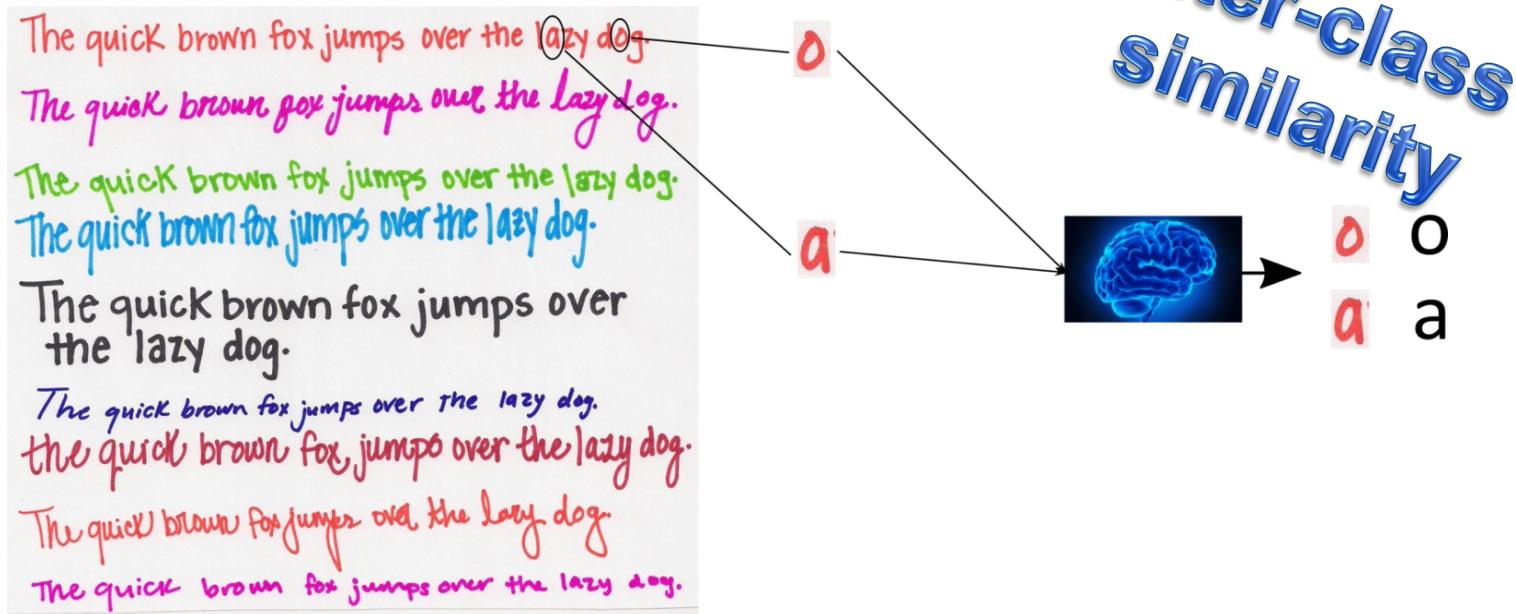
# Neural Networks

- Example: Handwritten characters classification
  - It is very simple for human beings to understand handwriting texts, even with the high inter-class similarity and high intra-class variance.



# Neural Networks

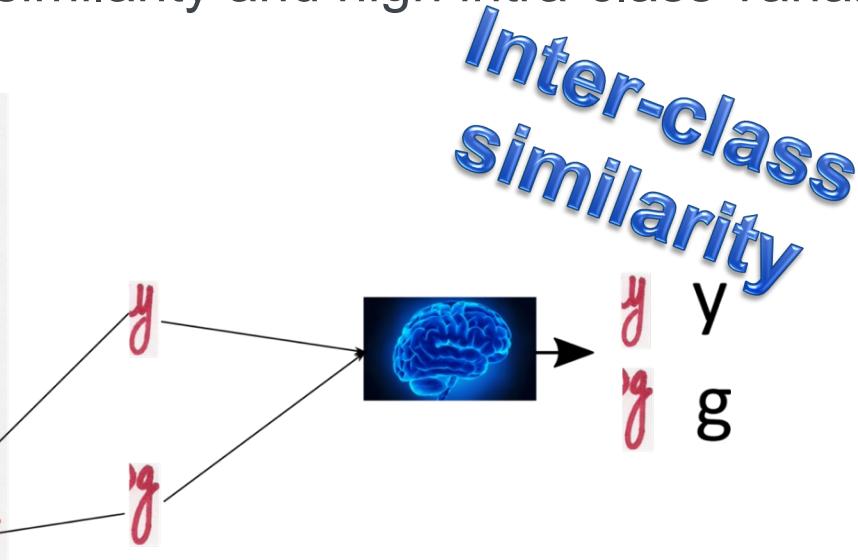
- Example: Handwritten characters classification
  - It is very simple for human beings to understand handwriting texts, even with the high inter-class similarity and high intra-class variance.



# Neural Networks

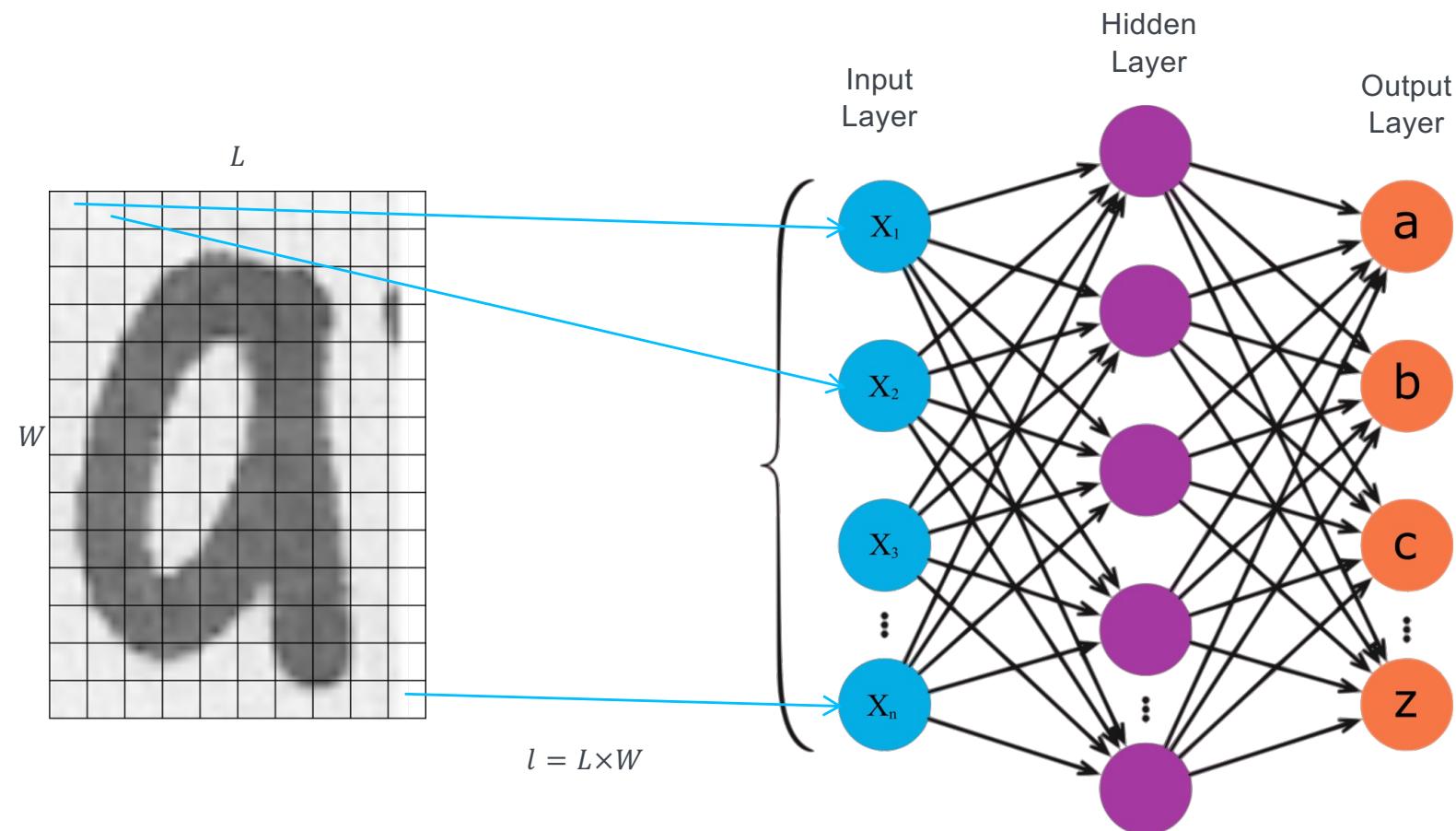
- Example: Handwritten characters classification
  - It is very simple for human beings to understand handwriting texts, even with the high inter-class similarity and high intra-class variance.

The quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.  
the quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.  
The quick brown fox jumps over the lazy dog.



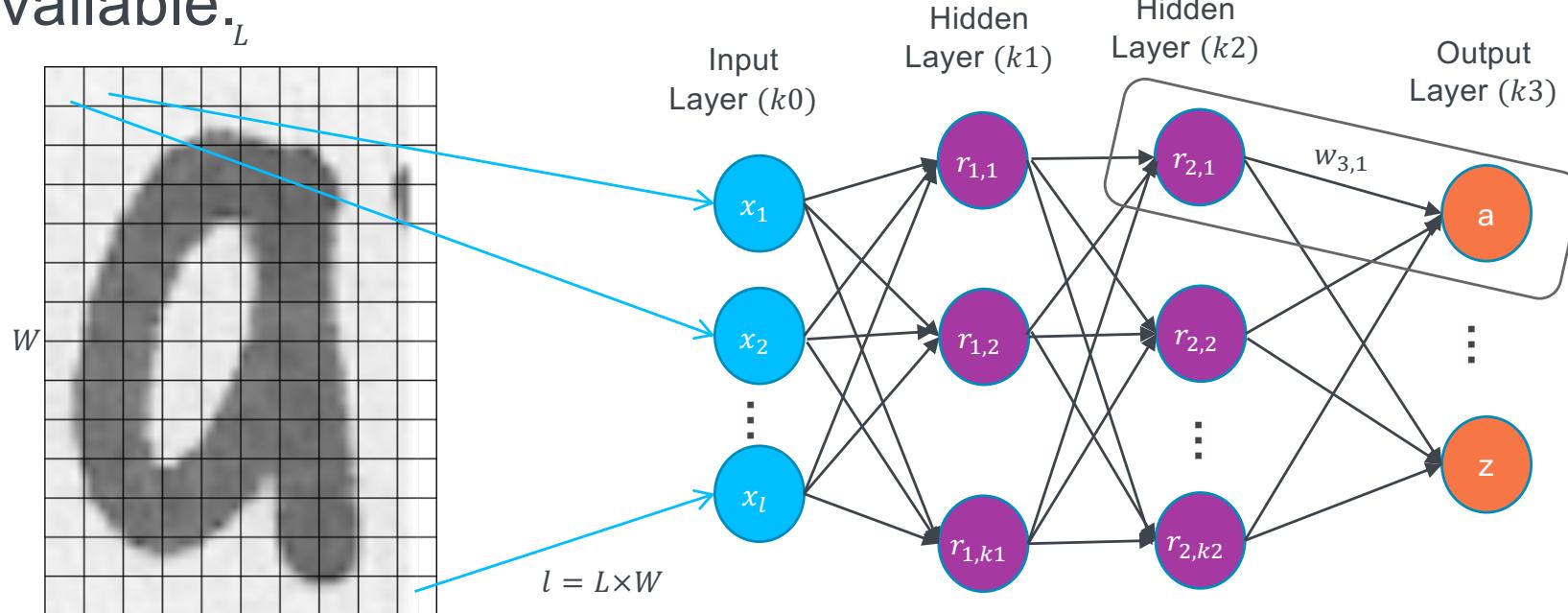
Inter-class  
similarity

# Neural Networks

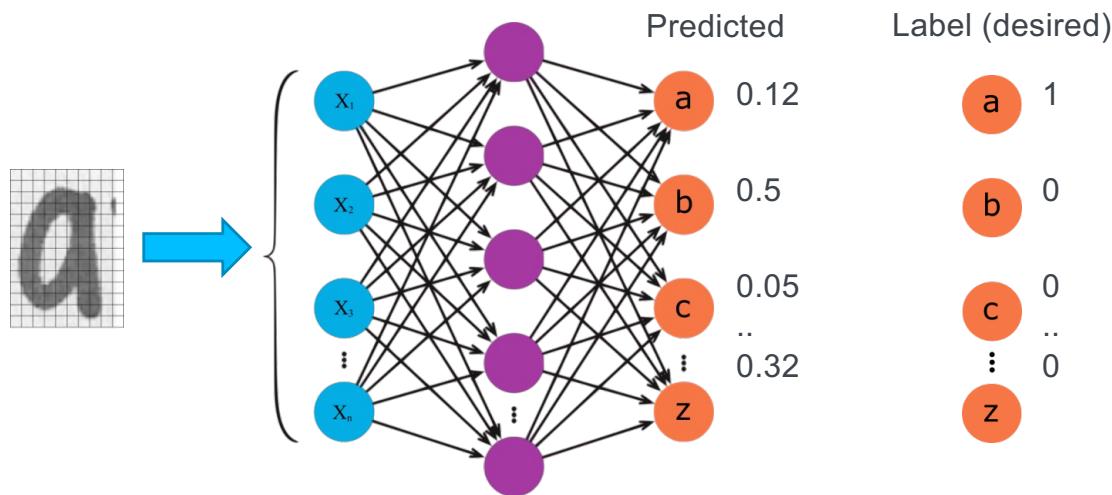


## How to build MLP?

- Given the architecture below, where all neurons employ the same activation function (sigmoid)
- We assume that  $N$  training samples  $(x_i, y_i), i = 1, \dots, N$  are available.



# Neural Networks



Loss function

# Predicting probability distributions

Output Type	Output Distribution	Output Layer	Cost Function
Binary	Bernoulli	Sigmoid	Binary cross-entropy
Discrete	Multinoulli	Softmax	Discrete cross-entropy
Continuous	Gaussian	Linear	Gaussian cross-entropy (MSE)
Continuous	Mixture of Gaussian	Mixture Density	Cross-entropy
Continuous	Arbitrary	GAN, VAE, FVBN	Various

# **6 Modeling a Feed-forward network predicting a Bernoulli distribution**

## Loss function with maximum likelihood

Given training data  $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^d, y_i \in \{0,1\}$

### Cross entropy

(as very often used for classification!)

$$L(\theta) = \mathbb{E}_{data}(-\log P(y|x)) = -\frac{1}{N} \sum_{i=1}^N \log P(y_i|x_i)$$

## Bernoulli distribution

Given training data  $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^d, y_i \in \{0,1\}$

Model  $z = w^T x + b$ , model  $\log P(y|x) \sim yz$

Using the sigmoid activation function, we get

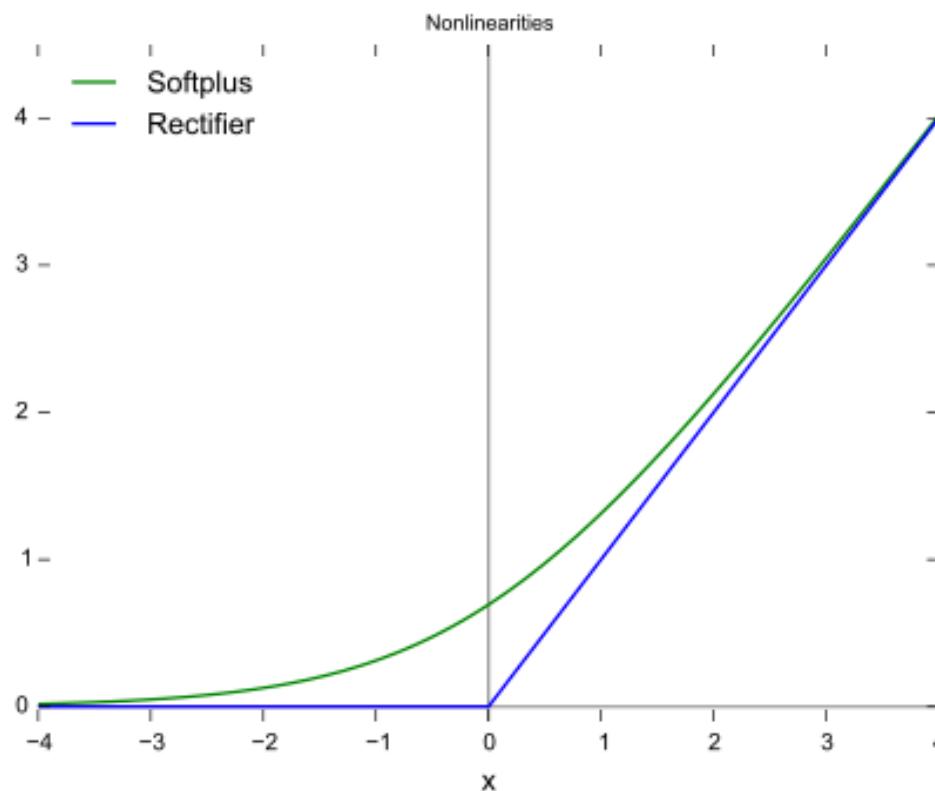
$$P(y|x) = \frac{e^{yz}}{\sum_{y'=0}^1 e^{y'z}} = \begin{cases} \text{Only large,} \\ \text{if } y = 1 \text{ and } z \text{ is large} \\ \text{Normalization to 1} \end{cases}$$
$$= \frac{e^{(2y-1)z}}{1 + e^{(2y-1)z}} = \sigma((2y - 1)z)$$

## Loss

$$L(\theta) = -\log P(y|x) = -\log \sigma((2y - 1)z) = \zeta((1 - 2y)z)$$

# Softplus function $\zeta$

$$\zeta(x) = \log(1 + e^x)$$



CC0, <https://en.wikipedia.org/w/index.php?curid=48817276>

$$\frac{d\zeta(x)}{dx} = \sigma(x)$$

$$\log \sigma(x) = -\zeta(-x)$$

We consider  $\zeta((1 - 2y)z)$ :

- When does  $\zeta$  saturate?
  - i.e. not provide guidance to the gradient?
- Only if  $(1 - 2y)z$  very negative, i.e.
  - $y=1$  and  $z$  very positive
  - $y=0$  and  $z$  very negative i.e. when prediction correct!

## Summary: Learning Bernoulli output

- Model:

$$\hat{f}(x) = \sigma(w^T x)$$

- Loss function: Softplus

$$L(\theta) = \zeta((1 - 2y)z)$$

Saturation when  
correct result  
achieved

- Solving with gradient descent;

$$\frac{d\zeta(x)}{dx} = \sigma'(x)$$

# **7 Modeling a Feed-forward network predicting a Multinoulli distribution**

## Comparing Bernoulli vs Multinoulli

Prediction for Bernoulli:

$$\hat{y} = P(y = 1|x)$$

With  $\hat{y}$  being the predicted value from our neural network for a given input vector  $x$

We require

$$\hat{y} \in [0,1]$$

Prediction for Multinoulli:

$$\hat{y}_i = P(y = i|x)$$

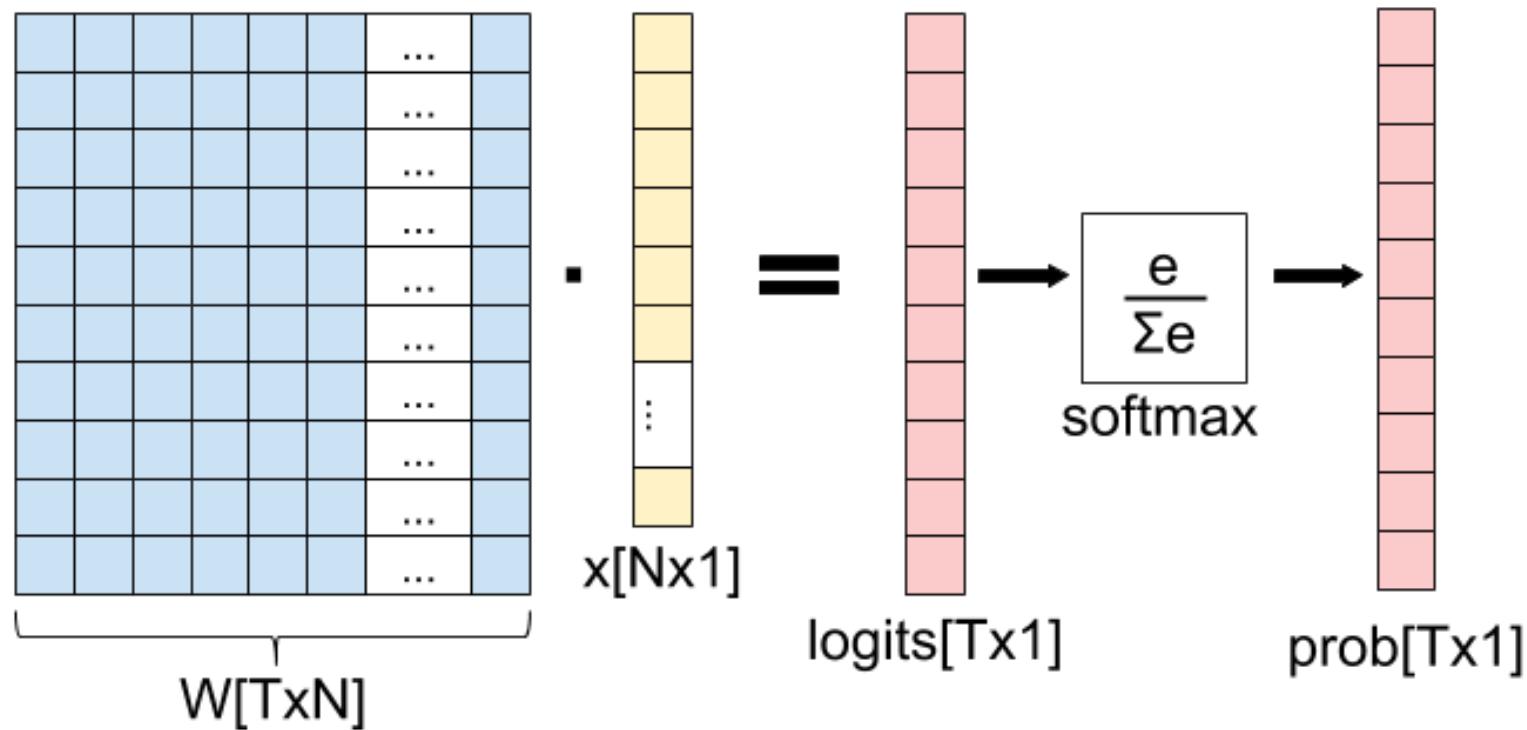
With  $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$  being the vector of predicted values from our neural network for  $n$  different classes

Now we require

$$\forall i: \hat{y}_i \in [0,1]$$

$$\sum_{i=1}^n \hat{y}_i = 1$$

## Multinoulli: Layers



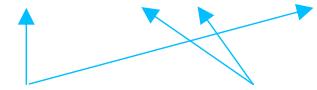
<https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/>

## Bernoulli vs Multinoulli: Layers

Prediction for Bernoulli:

$$\hat{y} = P(y = 1 | \mathbf{x})$$

We have defined

$$z = \mathbf{w}^T \mathbf{h} + b$$


scalar      vector

Prediction for Multinoulli:

$$\hat{y}_i = P(y = 1 | \mathbf{x})$$

Now we define

$$z = \mathbf{W}^T \mathbf{h} + \mathbf{b}$$

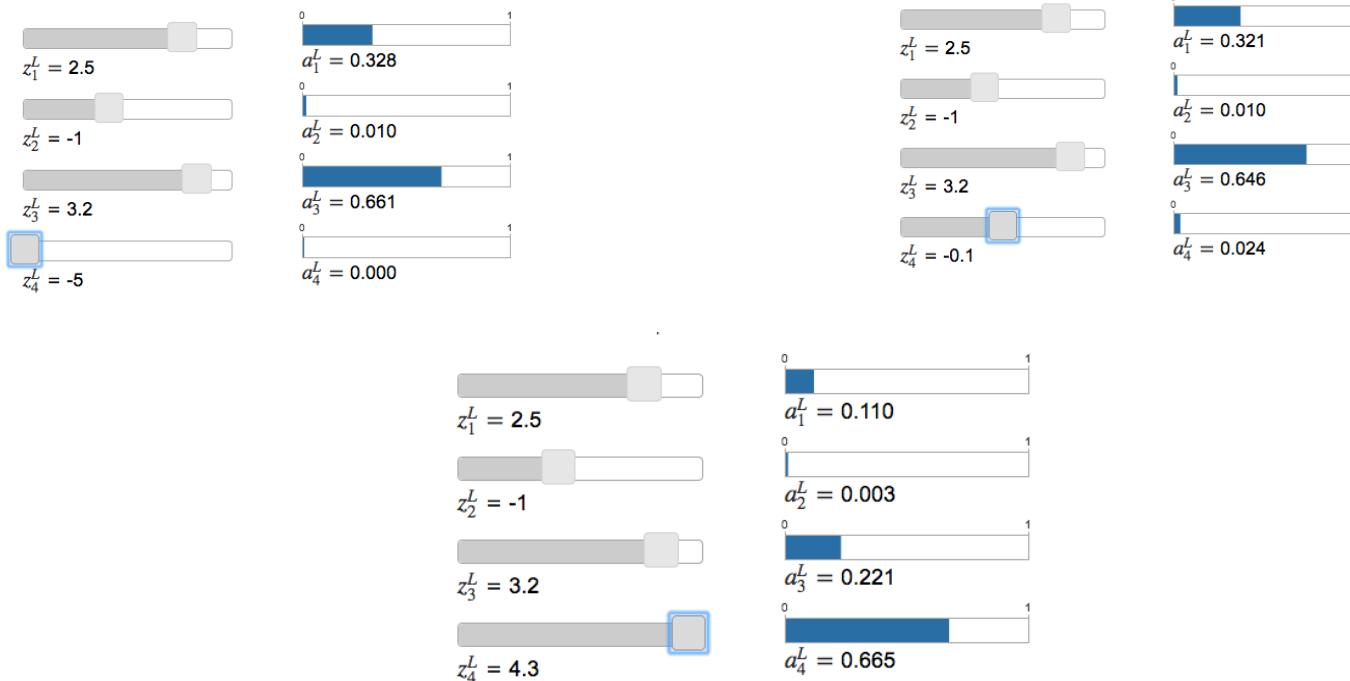

vector      matrix

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$


# Understanding the softmax

Animation at:

<http://neuralnetworksanddeeplearning.com/chap3.html#softmax>



## Multinoulli: Loss function

Cross entropy for a data record with  $y = (0 \dots 0_{i-1} 1_i 0_{i+1} \dots 0)$   
(often written in simplified manner as  $y = i$ )

$$\begin{aligned} & \text{data distribution} && \text{model distribution learned so far} \\ & - \sum_j 1_{j=i} \cdot \log \text{softmax}(\mathbf{z})_j = && \downarrow \\ & = -\log \text{softmax}(\mathbf{z})_i = && \\ & = -\log \frac{e^{z_i}}{\sum_j e^{z_j}} = && \\ & = -(\log e^{z_i} - \log \sum_j e^{z_j}) = && \\ & = -(z_i - \log \sum_j e^{z_j}) && \end{aligned}$$

## Stable softmax

```
In [146]: softmax([1, 2, 3])
Out[146]: array([ 0.09003057,  0.24472847,  0.66524096])
```

```
In [148]: softmax([1000, 2000, 3000])
Out[148]: array([ nan,  nan,  nan])
```

Remedy:

$$-\log \frac{e^{z_i}}{\sum_j e^{z_j}} = -\log \frac{e^{-c} e^{z_i}}{e^{-c} \sum_j e^{z_j}} = -\log \frac{e^{z_i - c}}{\sum_j e^{z_j - c}}$$

Pick:

$$c = \max_i z_i$$

## Summary on Multinoulli

- Softmax as soft generalization of argmax
- „Winner-takes-most“

## Side note

I liked the explanation by

<https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/>

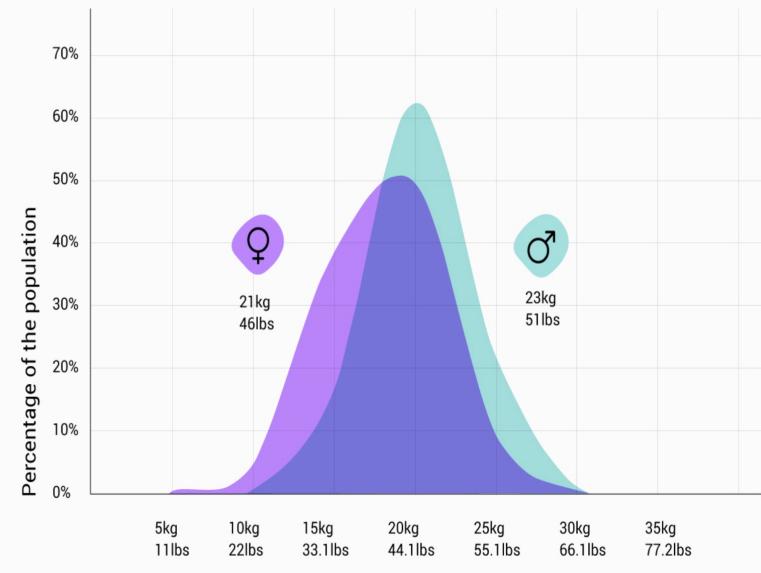
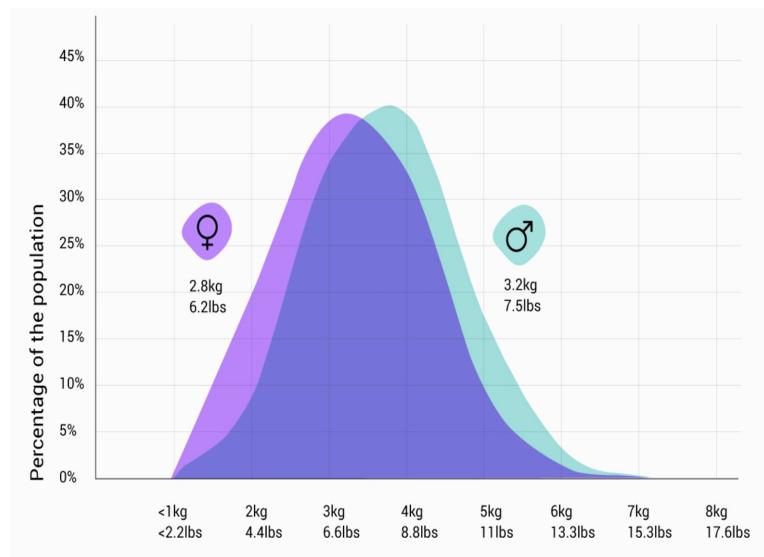
Better than the one in the neural network book

Combining linear output with multinoulli

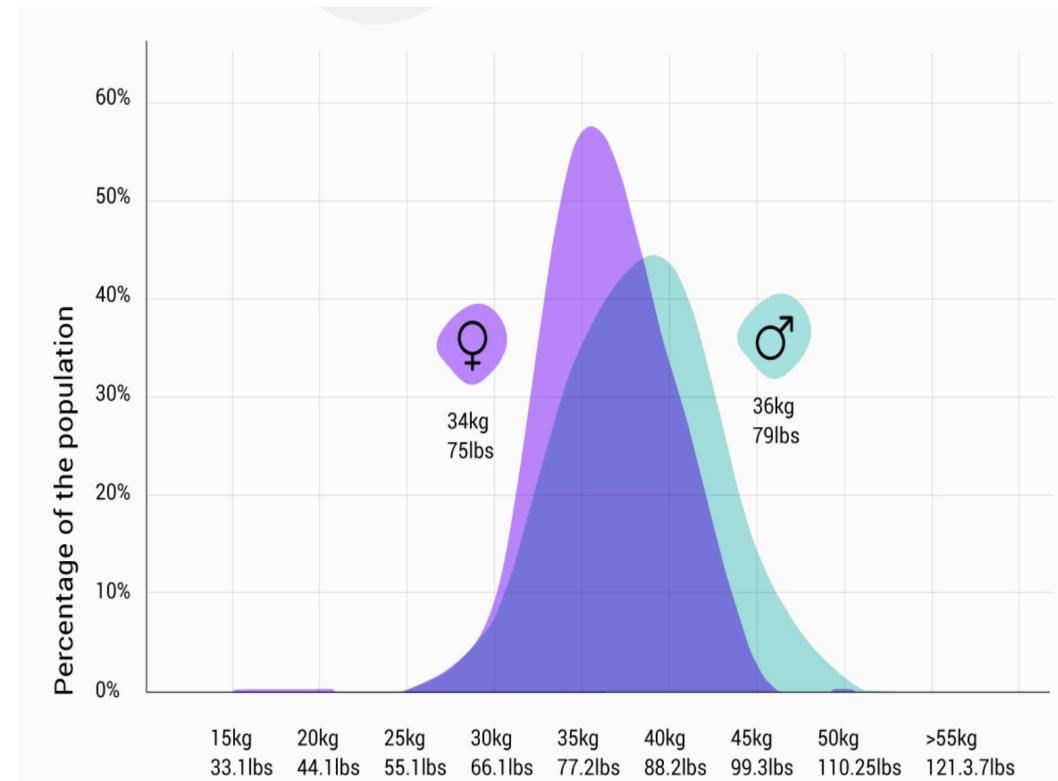
## 8 Modeling a Mixture Density Network

- (Bishop et al), p. 272
- (Goodfellow et al), chap. 6

# <https://www.11pets.com/en/news/siberian-husky-weight-chart>



<https://www.11pets.com/en/news/german-shepherd-weight-chart>



# Example

What do we want to achieve? Given an  $x$  we want to predict the most probable  $y$ .

Let's assume we have four dog breeds: from the tine Chihuahua, over Beagle, and Huskie, to the very large St. Bernard. Let's assume we have a set of dogs and we know some attributes of each dog like TypeOfFood ( $x_1$ ), Location ( $x_2$ ), age ( $x_3$ ), sex ( $x_4$ ) etc., and in the training set in particular we know their weight ( $y$ ), but we do \*not know\* their breed. Let's then predict the weight of a previously unseen dog, whose breed we do not know, but whose other attributes ( $x$  vector) we know.

We might do a linear regression to predict the weight which will vary from 1.8–2.7 kg (for the Chihuahua) to 54-82kg for the St. Bernard. This prediction will utterly fail, because for each of the four breeds we will have a very different mean of weight values.

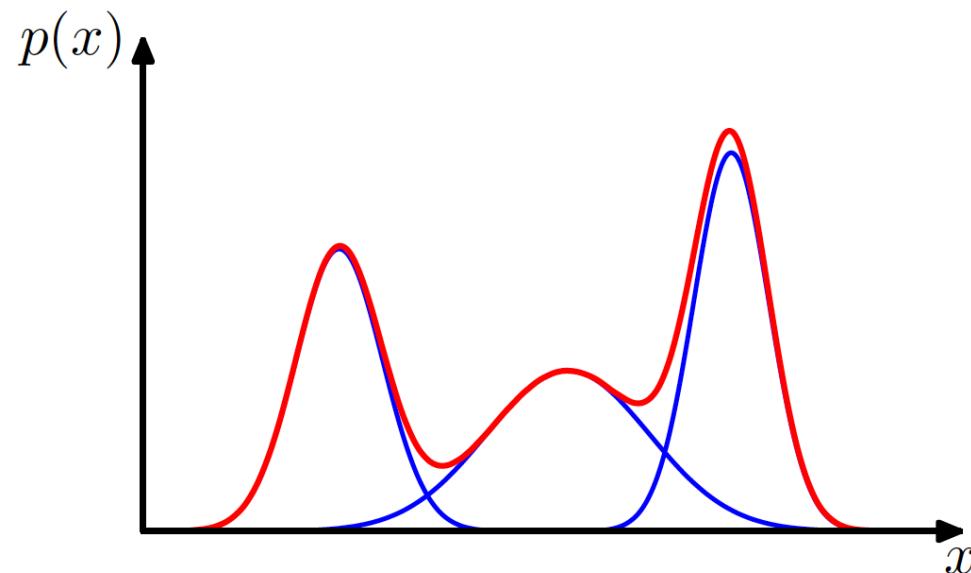
Hence, what we want to do is to consider the breed without knowing the breed. We want to discover in the data that there are four different clusters (around 2kg for the Chihuahua, 10kg for the Beagle, 20kg for the Husky, 70kg for the St. Bernard) and given the data (e.g. location and food) we might want to find that attributes like location tell us something about which breed is more probable than another breed and consider this for the prediction of the weight, where the weight for each cluster is approximately normal distributed around the center with larger variations along some axes (e.g. sex) and less variation along other axes.

The weight of each breed is approximately Gaussian distributed. Hence, trying to do both at the same time, i.e. finding the implicit cluster, its mean and its variation along the axes (e.g. St. Bernard's weight varies in the range of 64-82kg for males and 54 to 64kg for females, this is way more variation than one could find for the Chihuahua) allows us to do a joint prediction given some other attributes of the most probable weight value ( $y$ ).

Actual usage example would be in the area of prediction of vehicle movements, but there the assumptions are much more difficult than in the (unrealistic) dog example.

## Regression of a Gaussian Mixture

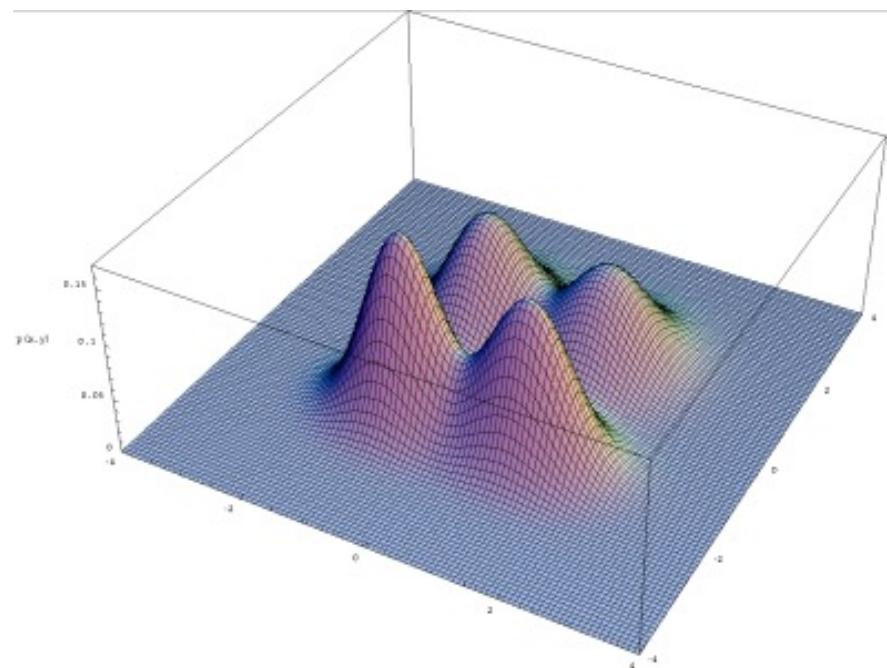
Univariate mixture distribution,  
showing bimodal distribution



(Bishop et al.), page 111

# Regression of a Gaussian Mixture

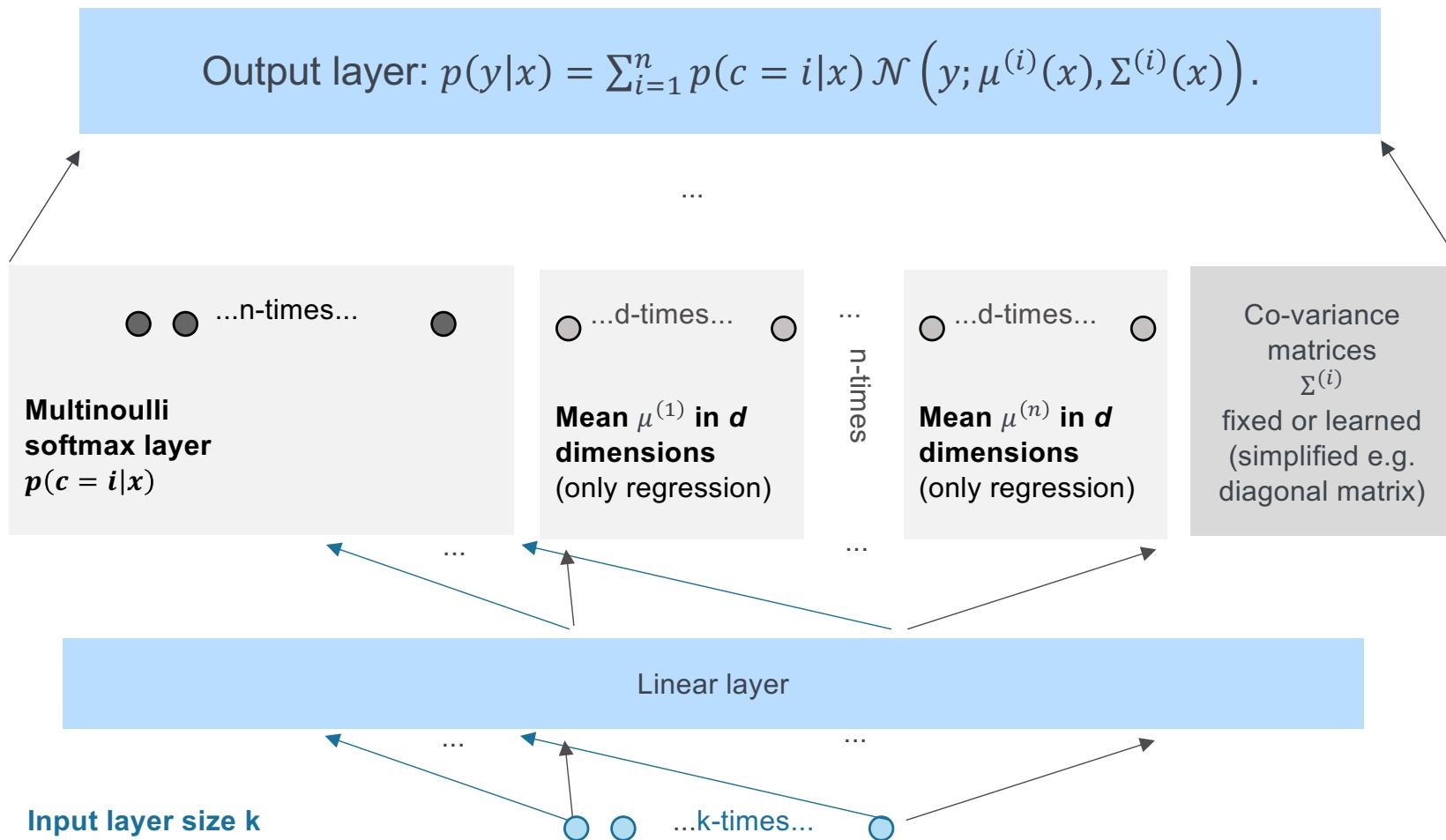
Multivariate (bi-variate) mixture distribution,  
showing four modes



By Jorgenumata (talk) (Uploads) (Own work) [Public domain], via Wikimedia Commons

<https://upload.wikimedia.org/wikipedia/commons/1/12/Bimodal-bivariate-small.png>

# Regression of a Gaussian mixture



## Predicting

$$\mathbb{E}[y|x] = \int y P(y|x) dy = \sum_{i=1}^n P(c=i) \mu^{(i)}(x)$$

# Gaussian mixture

- Different strategies for co-variance matrix
  - Constant
  - Diagonals (with partial influence according to category)
  - Full (with partial influence according to category)
- Loss function:
  - Mean Squared Error
    - Influences the means
    - Influences the choice of categories
      - make the dependence continuous
    - A mean should only be corrected if its category was influential

**Watch out in winter term 2025/2026!**

Seminar

**Machine Learning on Graphs**

by

Hernandez and team

Fachpraktikum

**Deep Learning Lab**

Nayyeri et al

**Watch out in winter term 2025+26!**

Lecture + Project

**Foundation Models / Deep Generative Models**

by

Mathias Niepert, Steffen Staab, Kai Arras

from BERT to ChatGPT, DALL-E, Llama etc.

**Limited participation**



Universität Stuttgart  
KI

# Thank you!



**Steffen Staab**

E-Mail [Steffen.staab@ki.uni-stuttgart.de](mailto:Steffen.staab@ki.uni-stuttgart.de)  
Telefon +49 (0) 711 685-88100  
[www.ki.uni-stuttgart.de/](http://www.ki.uni-stuttgart.de/)

Universität Stuttgart  
Analytic Computing, KI  
Universitätsstraße 32, 50569 Stuttgart