

دورهی آموزشی «علم داده»  
Data Science Course

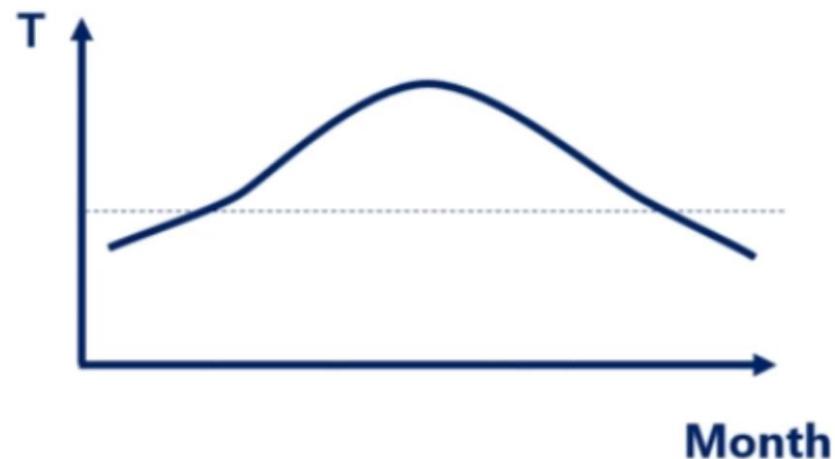


جلسه بیست و هشتم  
**مقدمه‌ای بر شبکه‌های عصبی عمیق** (همراه با معرفی کتاب)  
Deep Neural Networks

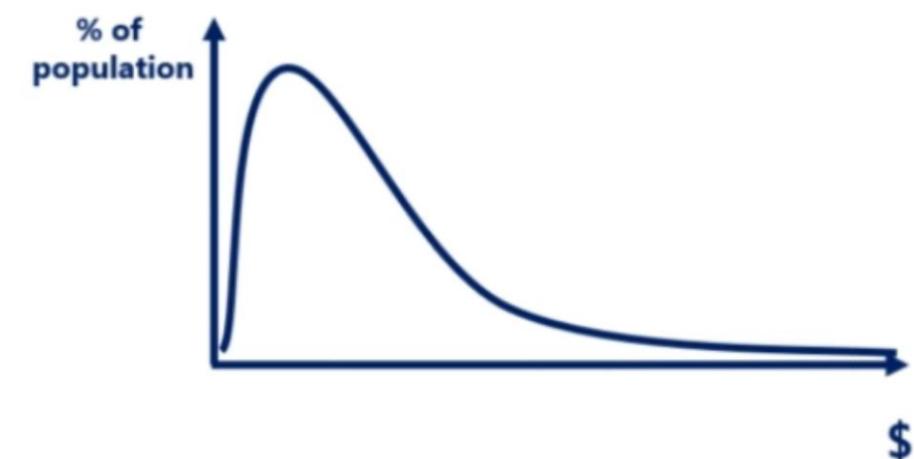
مدرس: محمد فزونی  
عضو هیئت علمی دانشگاه گنبدکاووس

# Real life is not linear

Temperature Northern hemisphere



Income distribution in the US



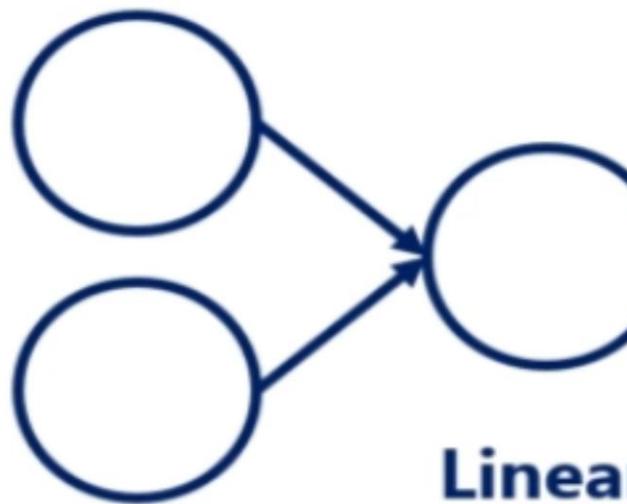
Because we want to be better forecasters, we need better models

**Linear**

**+**

**Non-linear**

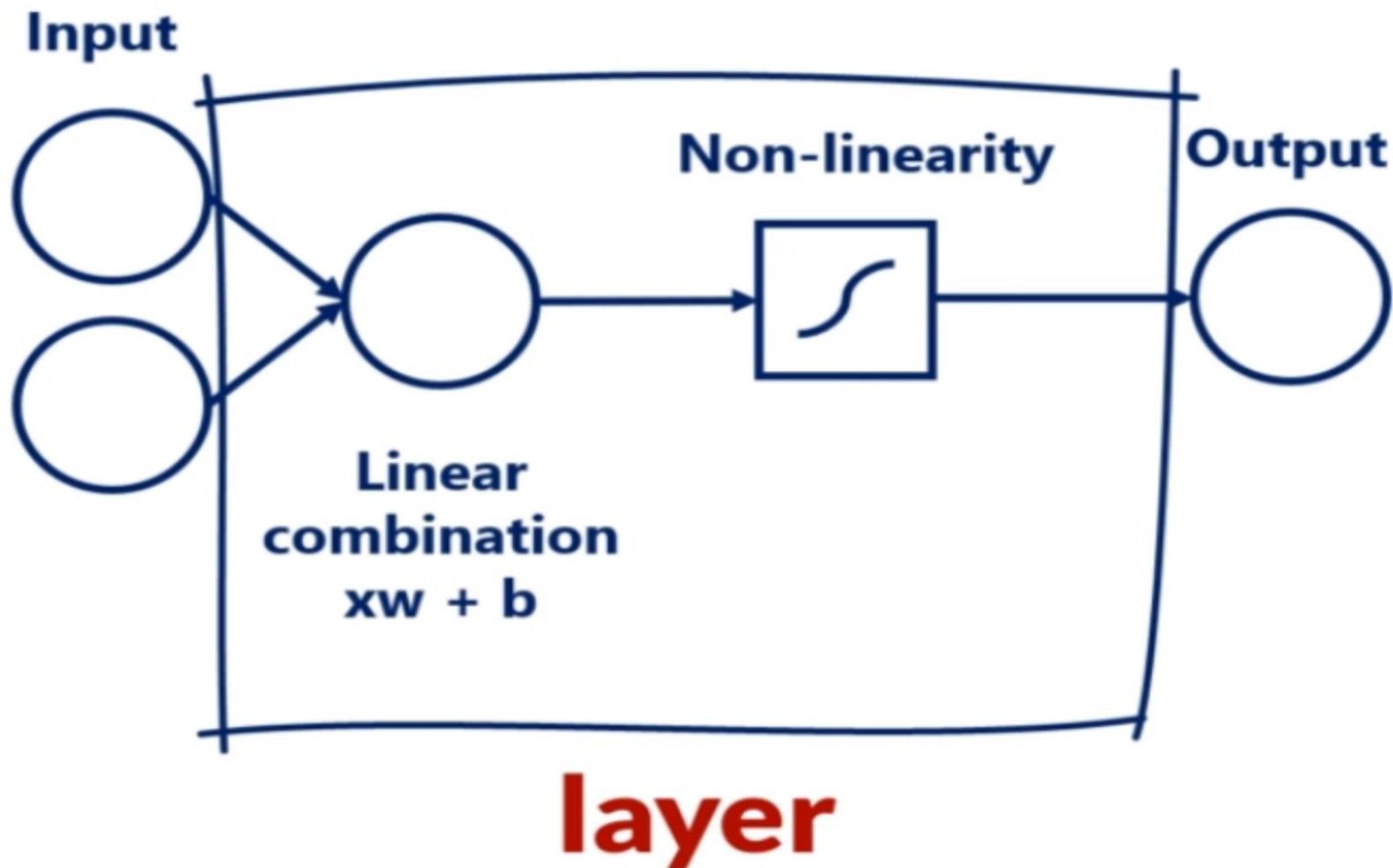
**Input**



**Non-linearity**

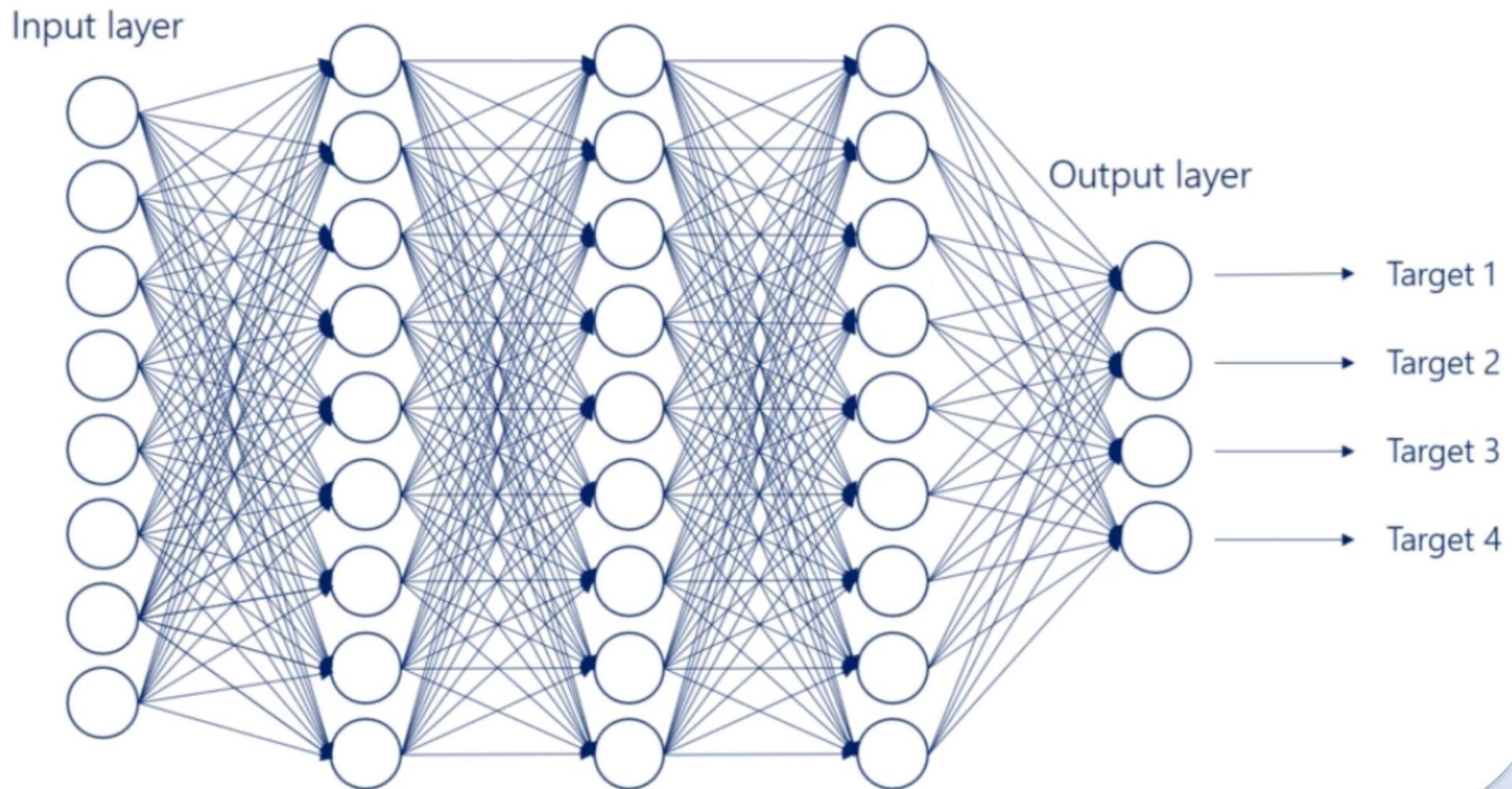
**Output**



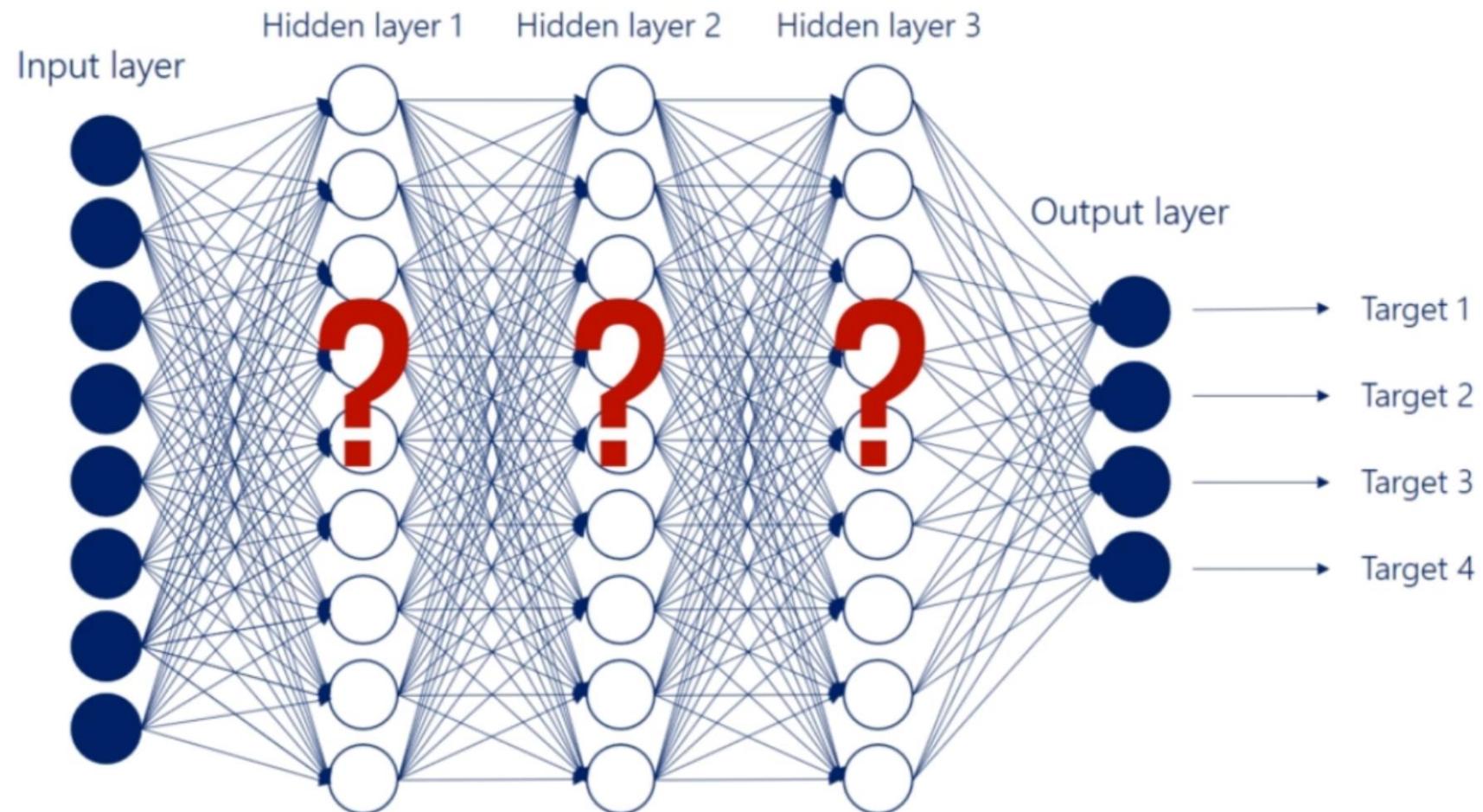


When we have more than one layer, we are talking about a **deep neural network**

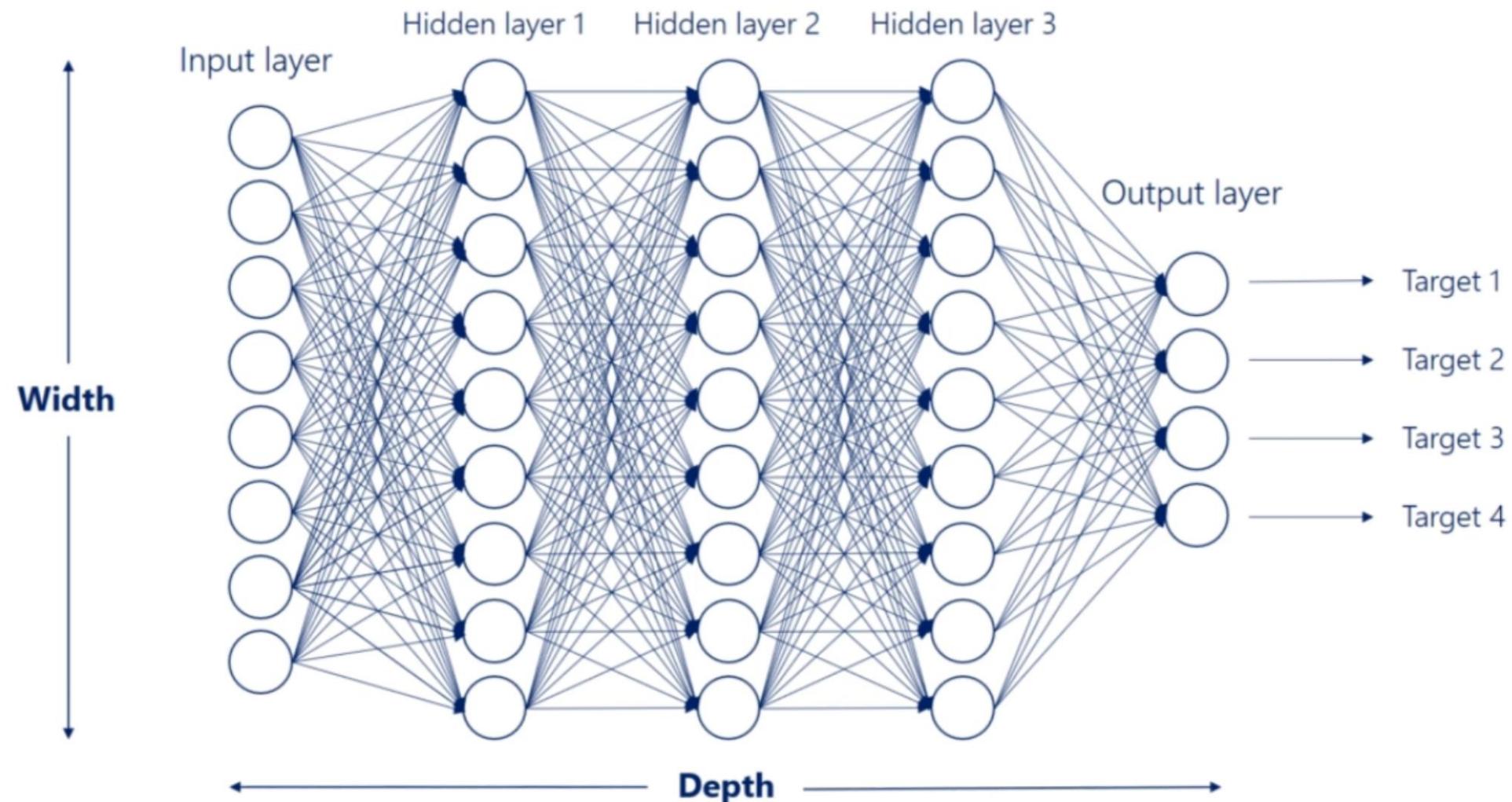
# What is a deep net?



# What is a deep net?



# A deep neural network



# Hyperparameters

vs

# Parameters

**Width**

**Depth**

**Learning rate (  $\eta$  )**

**Weights (  $w$  )**

**Biases (  $b$  )**

**Hyperparameters**  
**pre-set by us**

**vs**

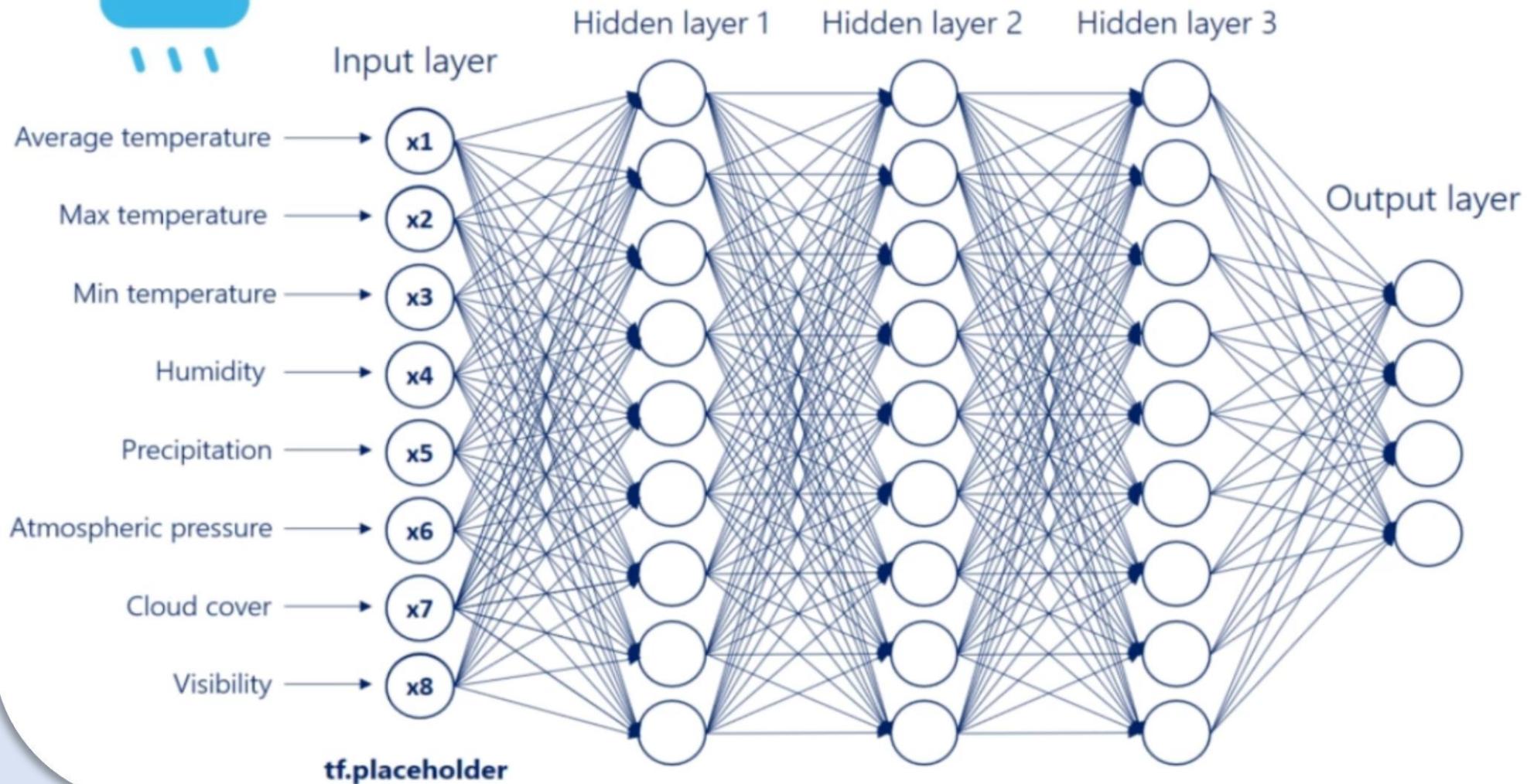
**Parameters**  
**found by optimizing**

**Width**  
**Depth**  
**Learning rate (  $\eta$  )**

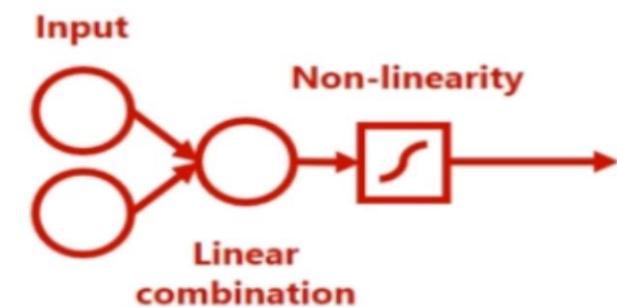
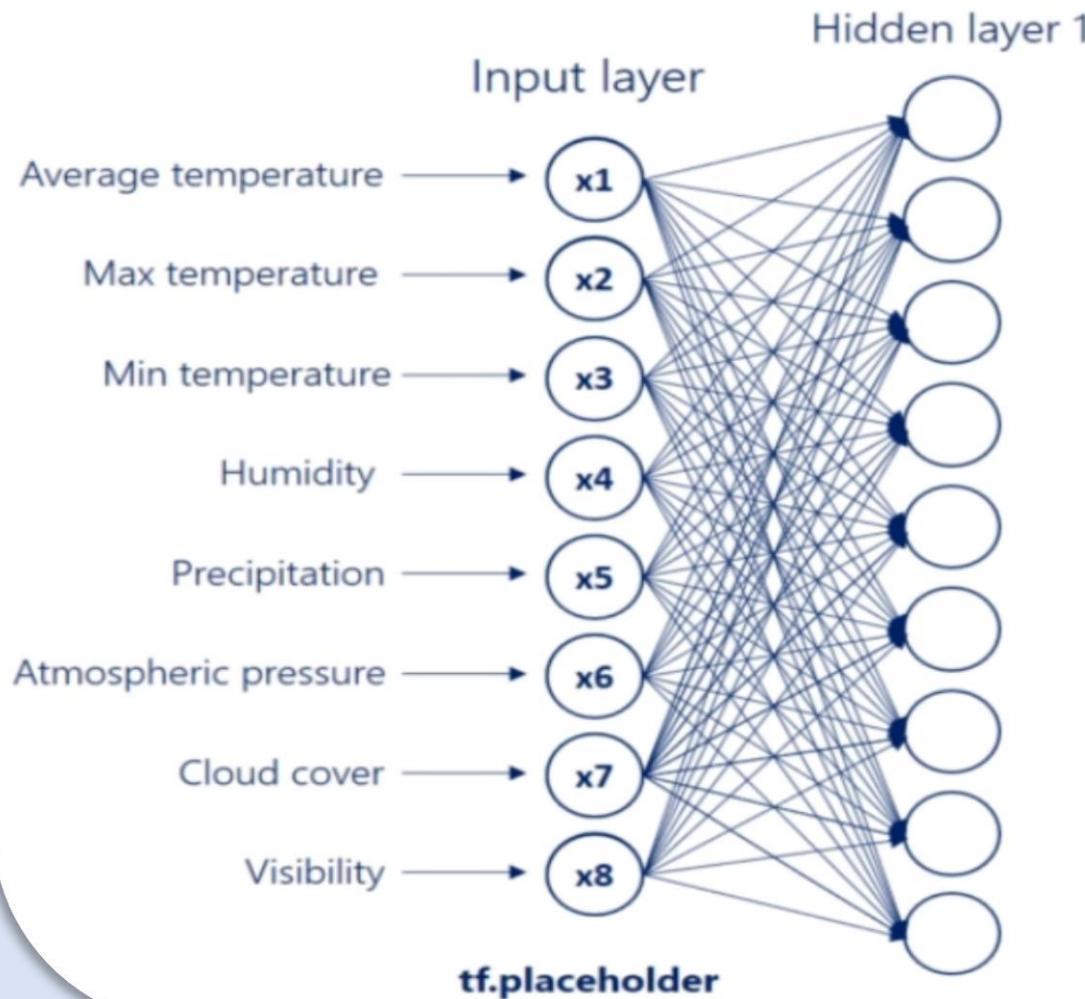
**Weights (  $w$  )**  
**Biases (  $b$  )**



# An illustration of deep nets

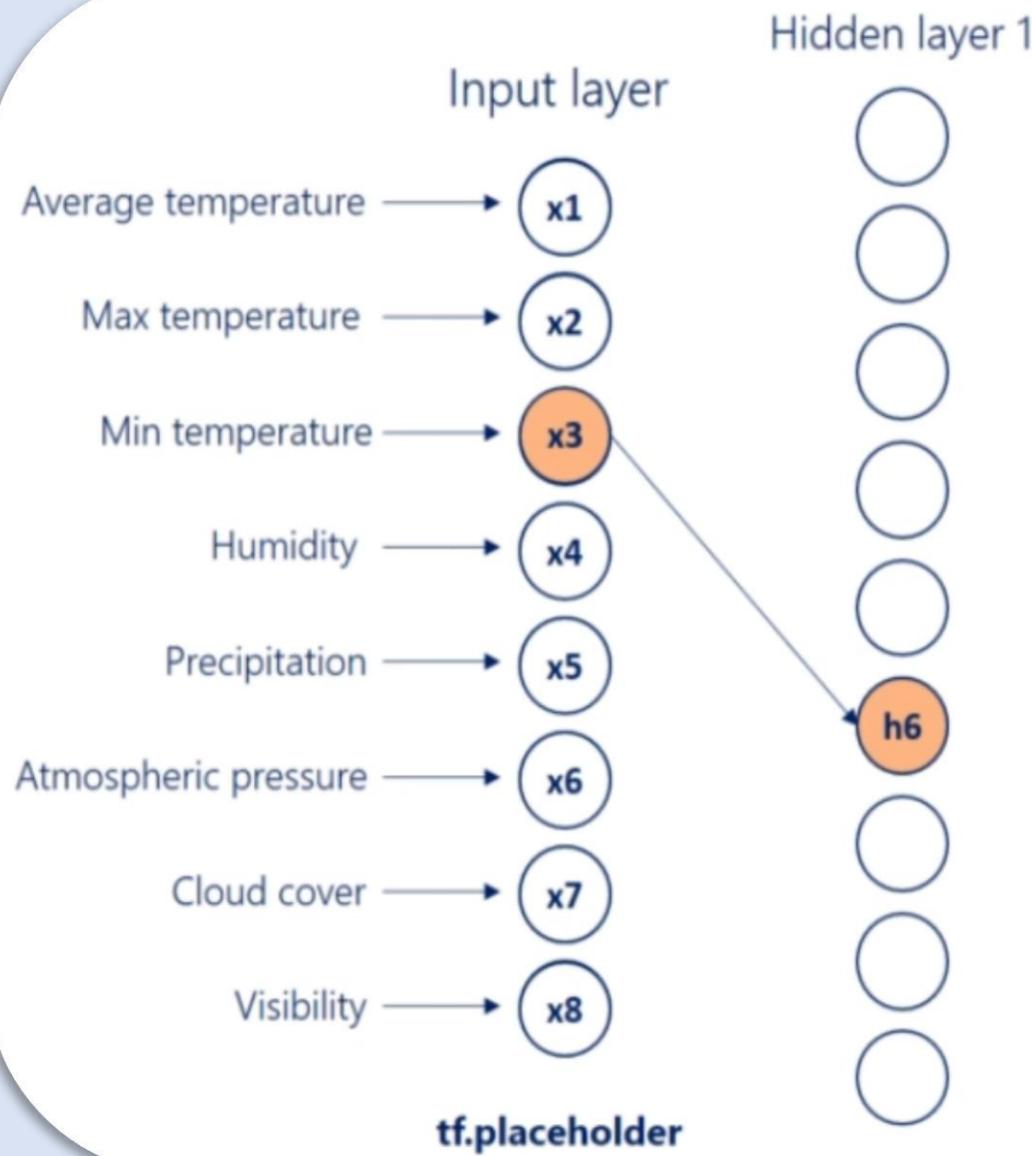


# An illustration of deep nets



$$X * W$$

$$1 \times 8 * 8 \times 9 = 1 \times 9$$

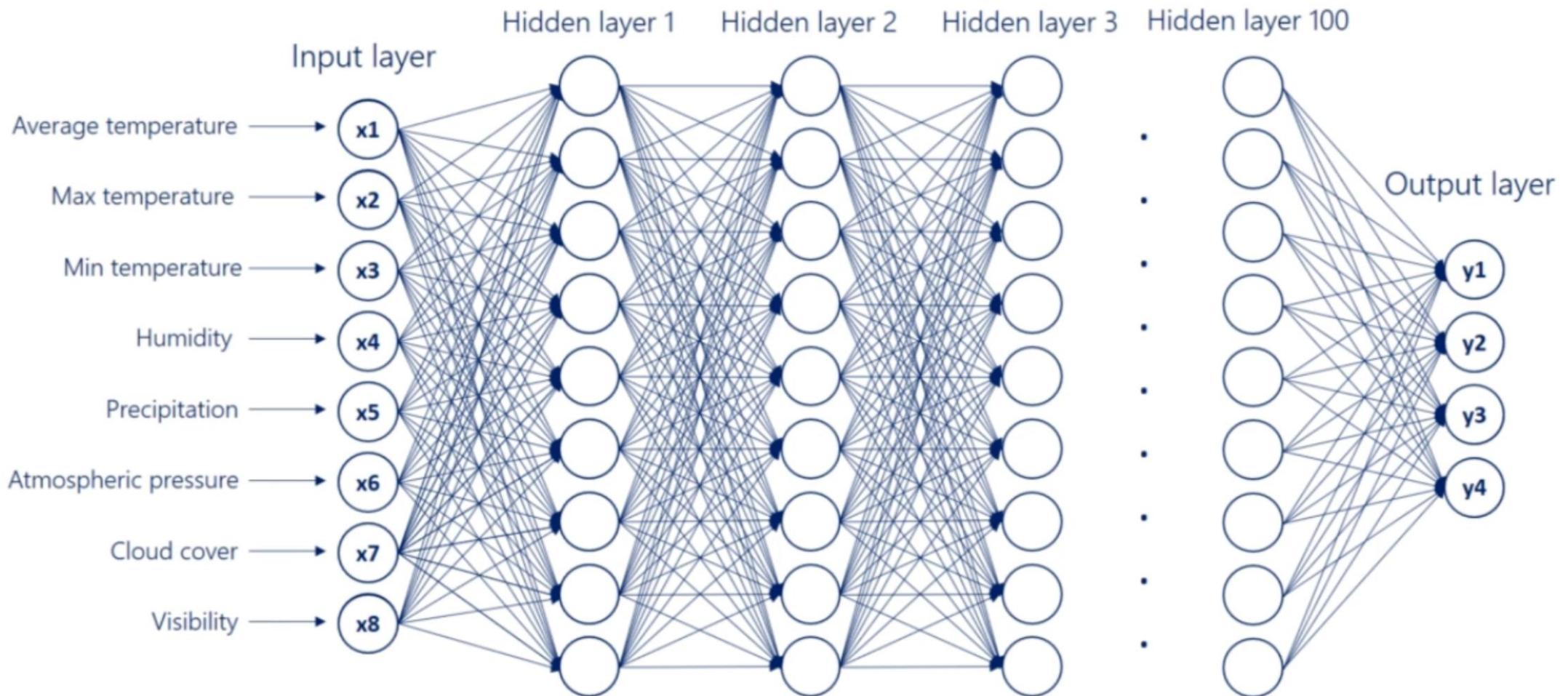


**W<sub>36</sub>**

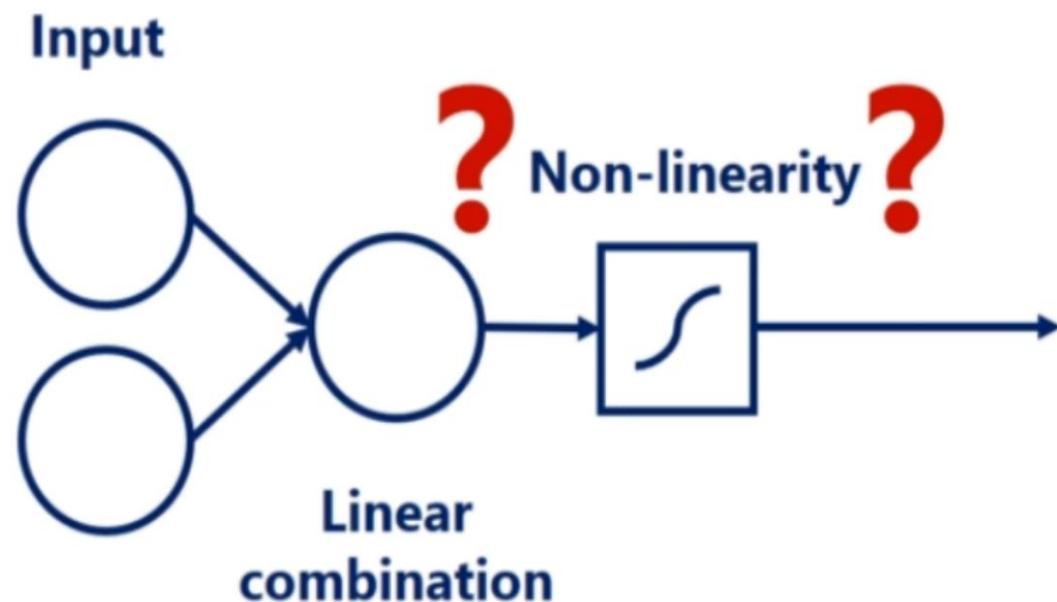
Input unit      Hidden unit

```
graph LR; subgraph W36 [W36]; direction TB; I[Input unit] --> W[W36]; H[Hidden unit] --> W; end;
```

# An illustration of deep nets



# Why do we need non-linearities?

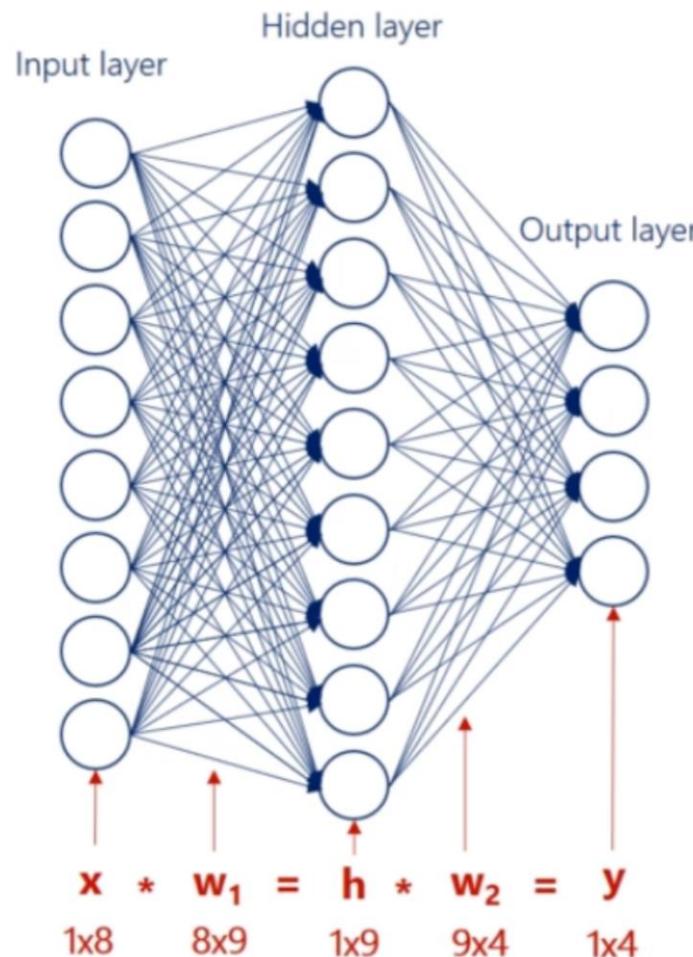


*"Non-linearities are needed so we can break the linearity and represent more complicated relationships"*

**We cannot stack layers when we  
only have linear relationships**

# Why do we need non-linearities?

Imagine a network with no **non-linearities**, just linear combinations



$$\mathbf{h} = \mathbf{x} * \mathbf{w}_1$$
$$\mathbf{y} = \mathbf{h} * \mathbf{w}_2$$

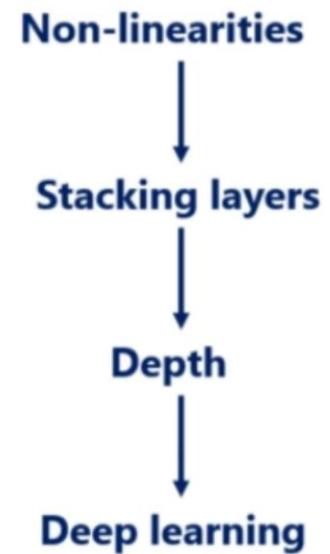
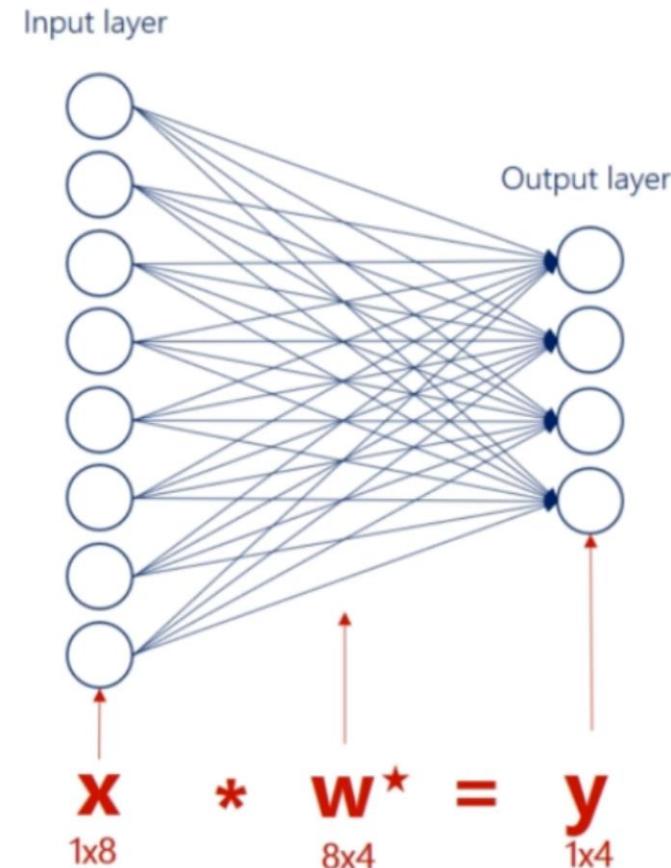
$$\begin{aligned} h &= x * w_1 \\ y &= h * w_2 \end{aligned}$$

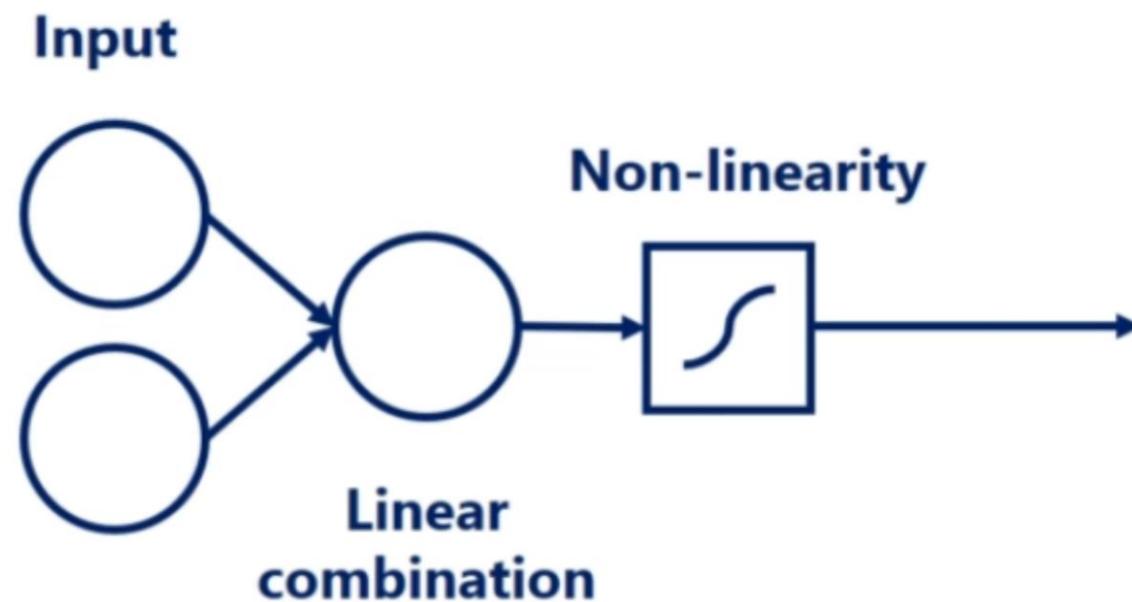
$$\begin{aligned} h &= x * w_1 \\ y &= x * w_1 * w_2 \end{aligned}$$

$$\begin{aligned} h &= x * w_1 \\ y &= x * w_1 * w_2 = \\ &\quad \boxed{\begin{matrix} 8 \times 9 & 9 \times 4 \end{matrix}} \\ &= x * w^* \\ &\quad \quad \quad 8 \times 4 \end{aligned}$$

# Why do we need non-linearities?

Two consecutive linear transformations are equivalent to a single one.





**In order to have deep nets and find complex relationships through arbitrary functions, we need non-linearities.**

یک مثال ساده از عملکرد جز غیر-خطی

Activation Function  
Non-linearity  
Transfer Function

**INPUT**



**ACTIVATION FUNCTION**



**Linear**



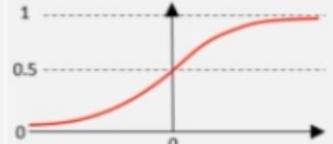
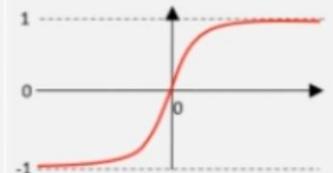
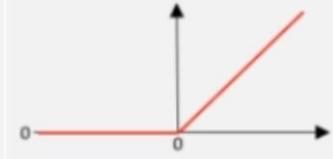
**OUTPUT**



**Non-linear**



## Common activation functions

Name	Formula	Derivative	Graph	Range
<b>sigmoid (logistic function)</b>	$\sigma(a) = \frac{1}{1+e^{-a}}$	$\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a))$		(0,1)
<b>TanH (hyperbolic tangent)</b>	$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$	$\frac{\partial \tanh(a)}{\partial a} = \frac{4}{(e^a + e^{-a})^2}$		(-1,1)
<b>ReLu (rectified linear unit)</b>	$\text{relu}(a) = \max(0,a)$	$\frac{\partial \text{relu}(a)}{\partial a} = \begin{cases} 0, & \text{if } a \leq 0 \\ 1, & \text{if } a > 0 \end{cases}$		(0,∞)
<b>softmax</b>	$\sigma_i(\mathbf{a}) = \frac{e^{a_i}}{\sum_j e^{a_j}}$	$\frac{\partial \sigma_i(\mathbf{a})}{\partial a_j} = \sigma_i(\mathbf{a}) (\delta_{ij} - \sigma_j(\mathbf{a}))$ Where $\delta_{ij}$ is 1 if $i=j$ , 0 otherwise	different every time	(0,1)

# یک مثال ساده از سافت‌مکس

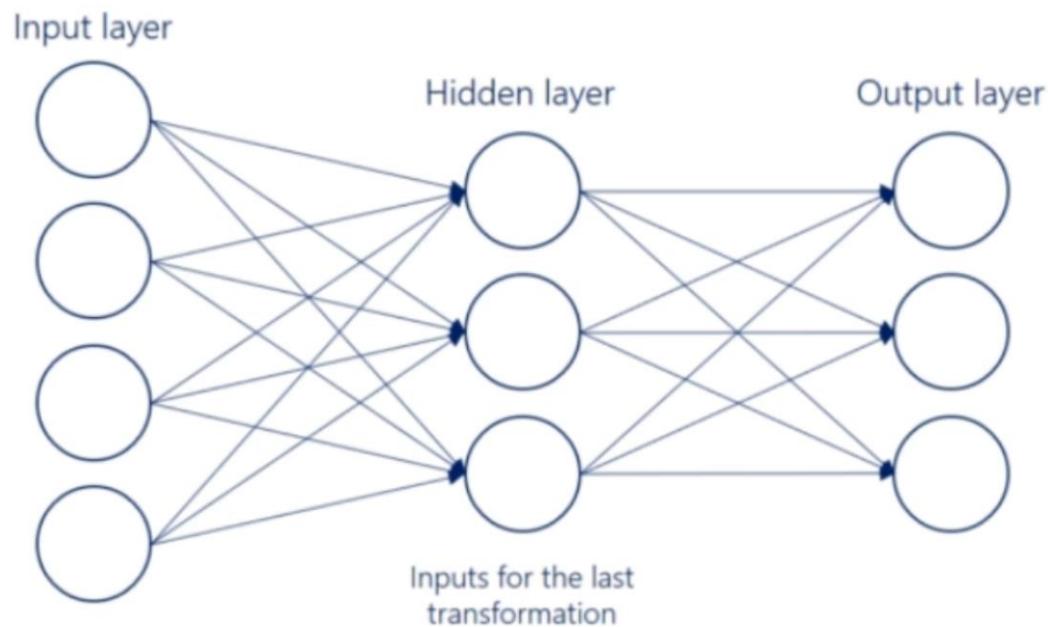
## Softmax

$$a = xw + b$$

Linear combination

$$y = \text{softmax}(\mathbf{a})$$

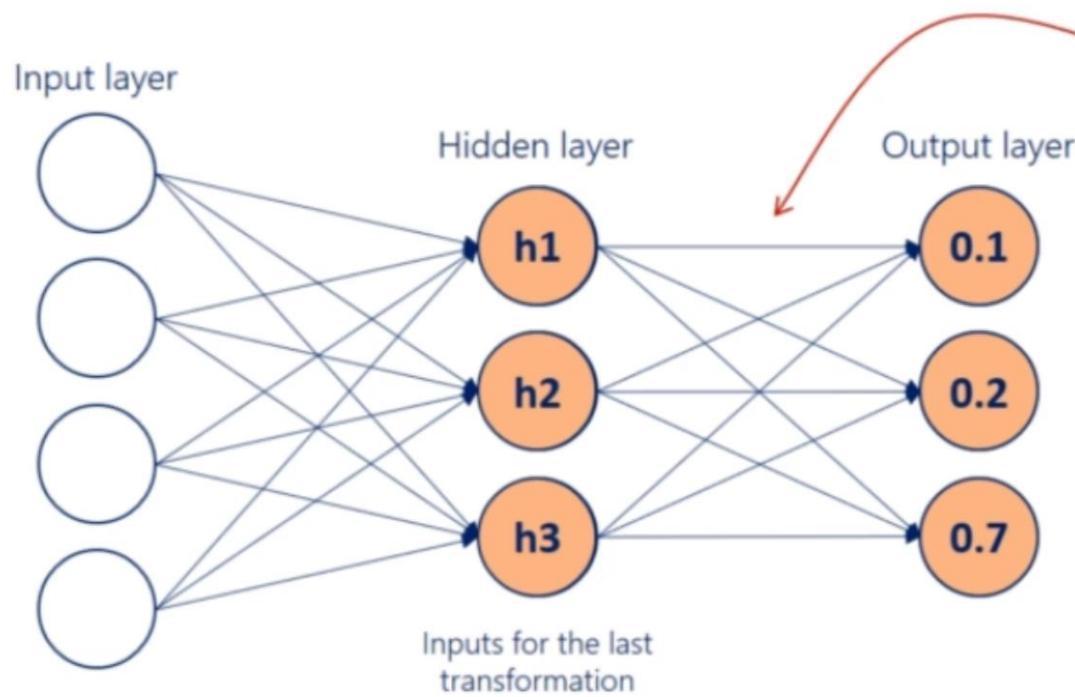
Activation



## Softmax

$$a = xw + b$$

Linear combination



$$y = \text{softmax}(\mathbf{a})$$

Activation

$$\mathbf{a}_h = hw + b$$

$$\mathbf{a} = [-0.21, 0.47, 1.72]$$

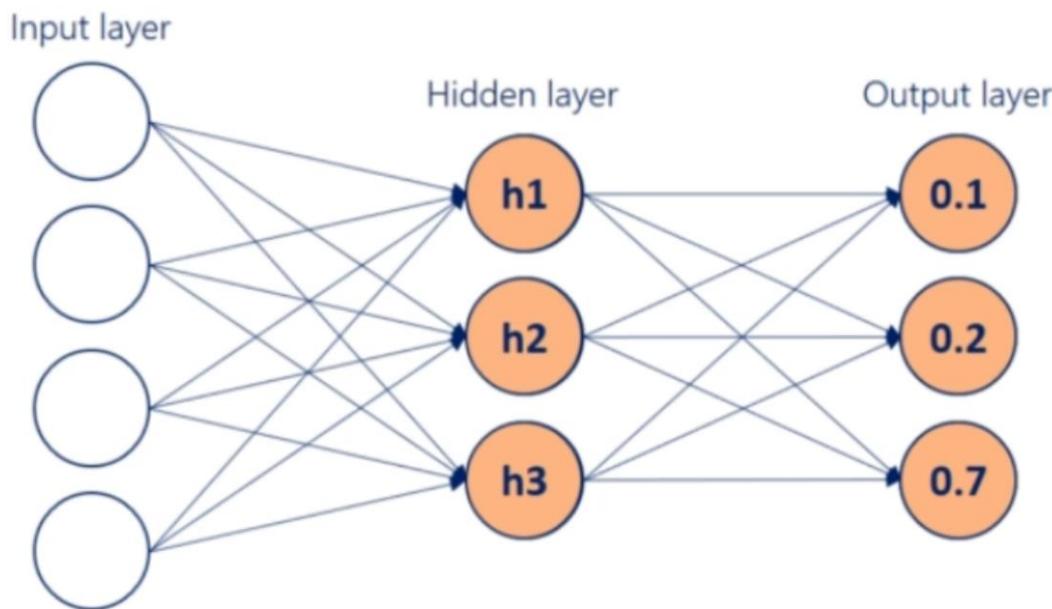
$$\text{softmax}(\mathbf{a}) = \frac{e^{a_i}}{\sum_j e^{a_j}}$$

$$\sum_j e^{a_j} = e^{-0.21} + e^{0.47} + e^{1.72} = 8$$

$$\text{softmax}(\mathbf{a}) = \left[ \frac{e^{-0.21}}{8}, \frac{e^{0.47}}{8}, \frac{e^{1.72}}{8} \right]$$

$$y = [0.1, 0.2, 0.7]$$

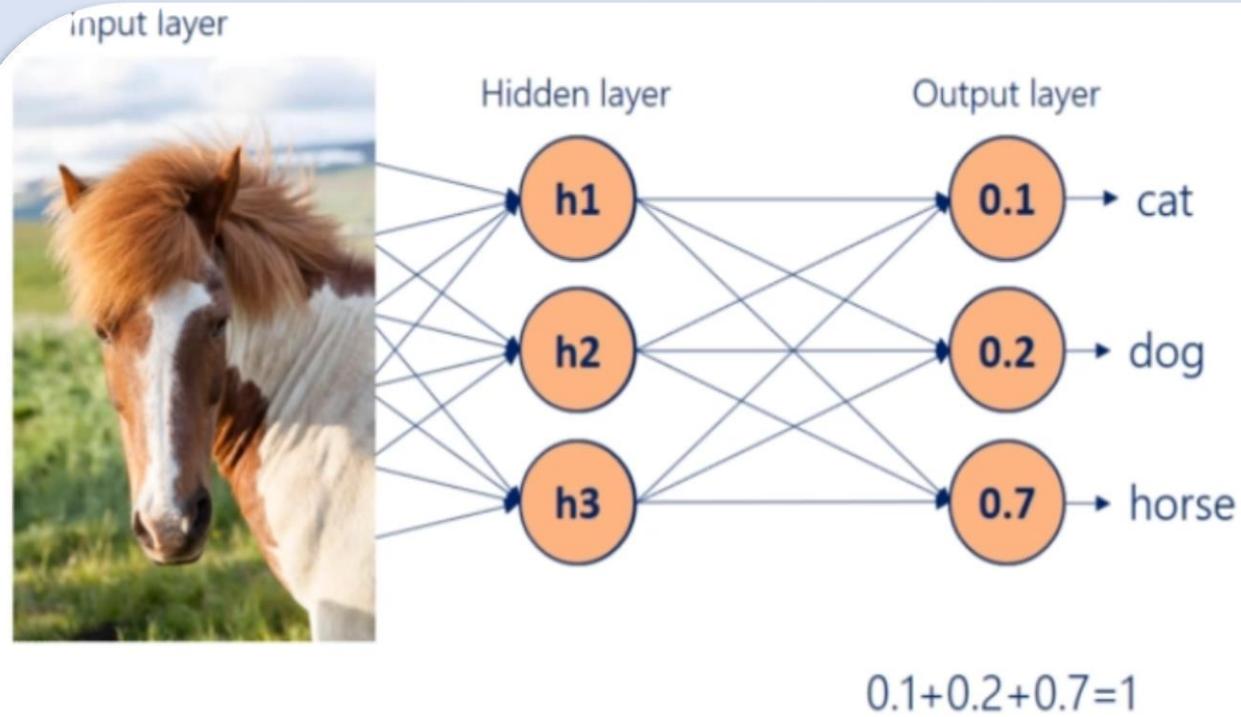
## Softmax



1. Range: (0,1)
2. They always sum up to 1

## PROBABILITIES

The softmax transformation transforms a bunch of arbitrarily large or small numbers into a valid probability distribution

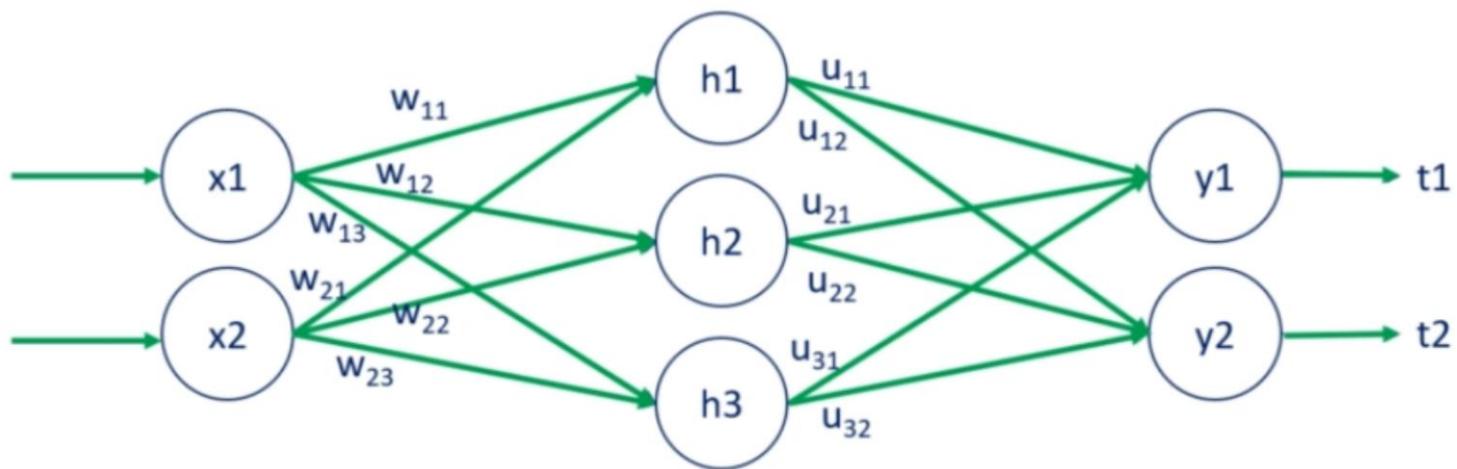


# PROBABILITIES

The softmax transformation transforms a bunch of arbitrarily large or small numbers into a valid probability distribution

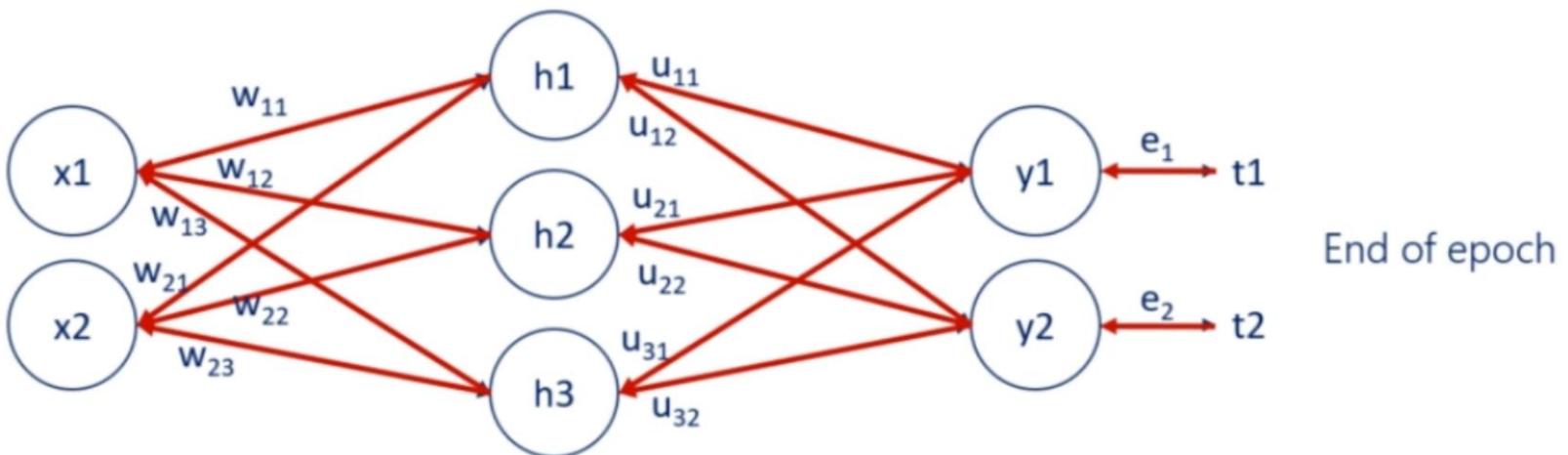
Softmax is often used as the activation of the output layer in classification problems

## Forward propagation



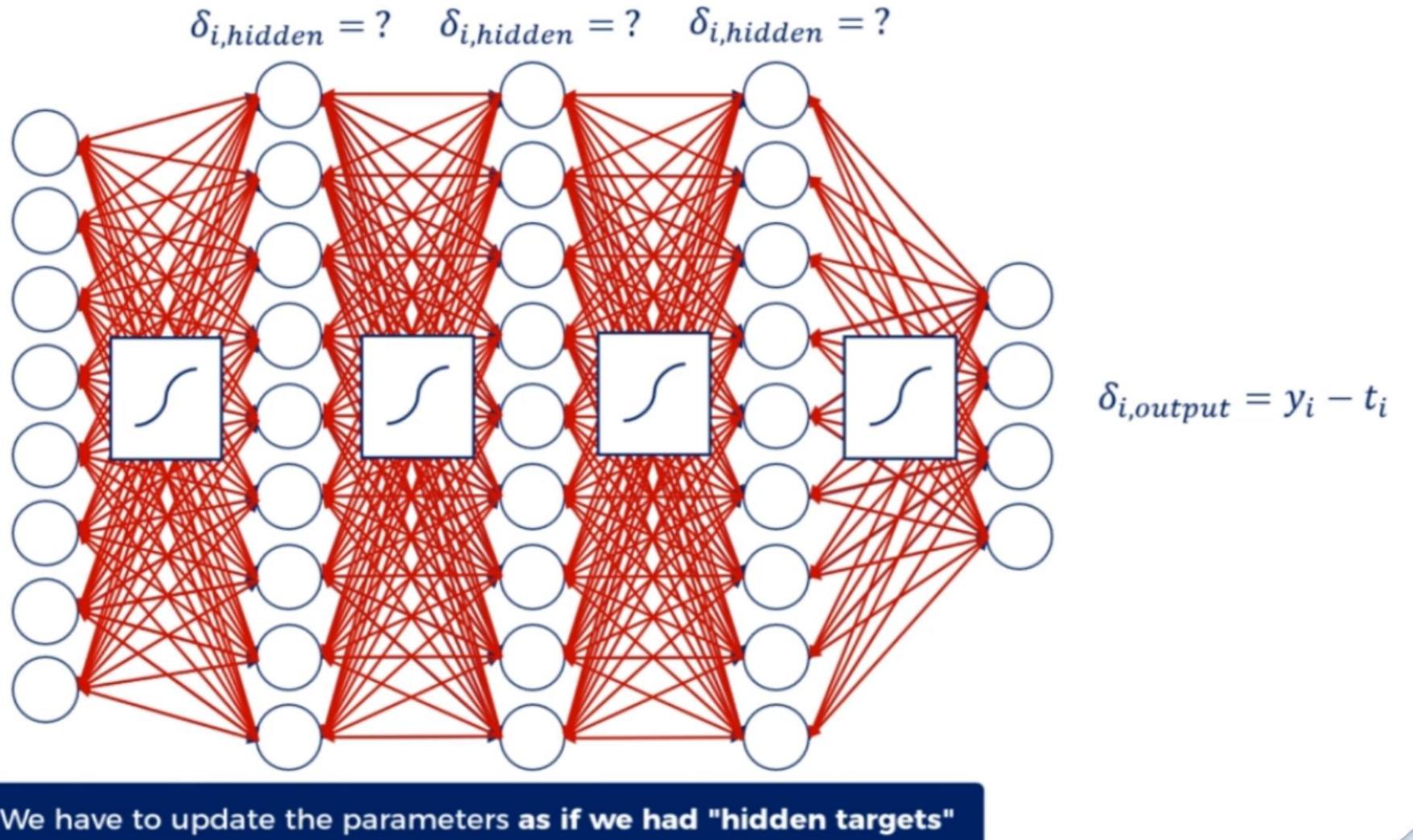
**Forward propagation** is the process of pushing inputs through the net

# Backpropagation

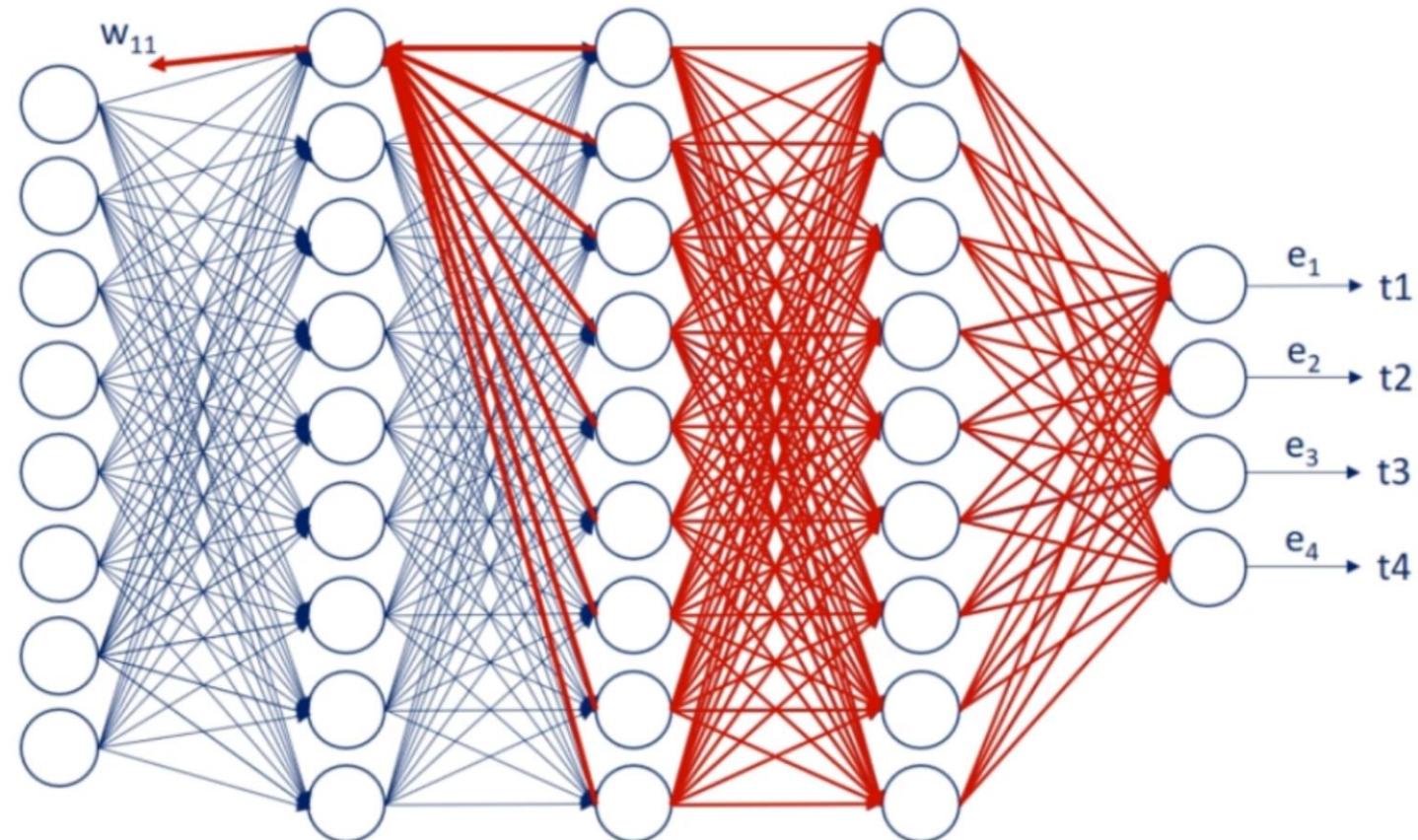


At the end of each epoch, we **backpropagate** and change each parameter accordingly

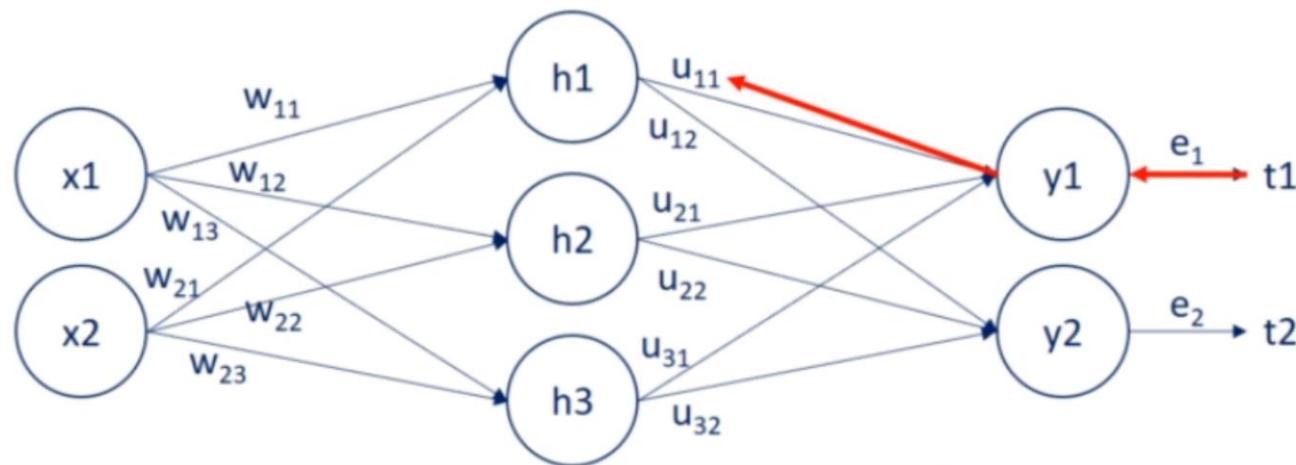
# Backpropagation



# Backpropagation



# Backpropagation

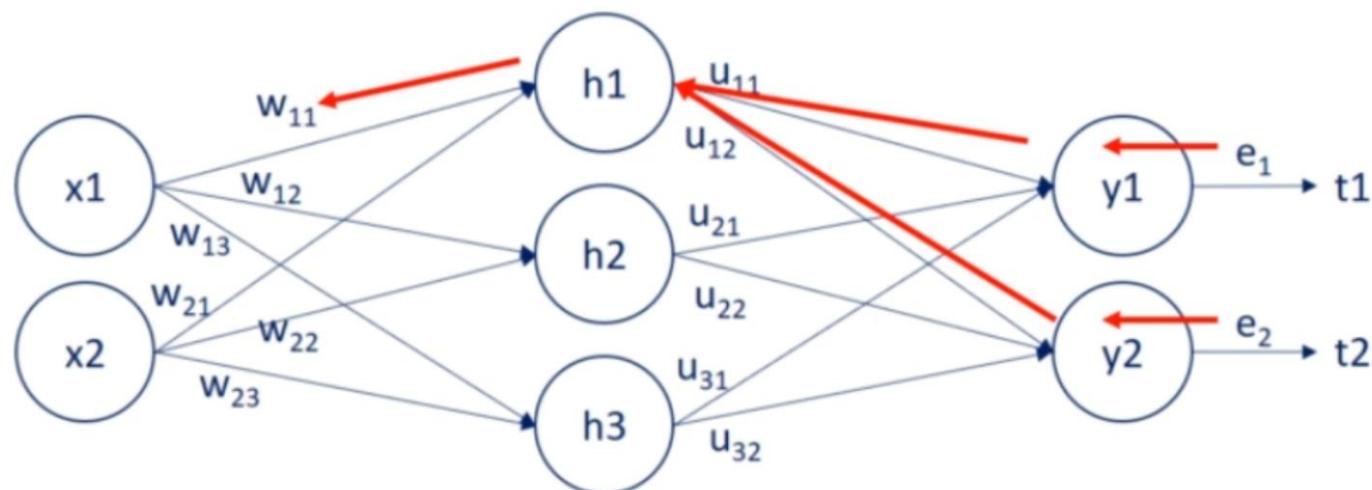


**$U_{11}$  contributes to a single error  $e_1$**

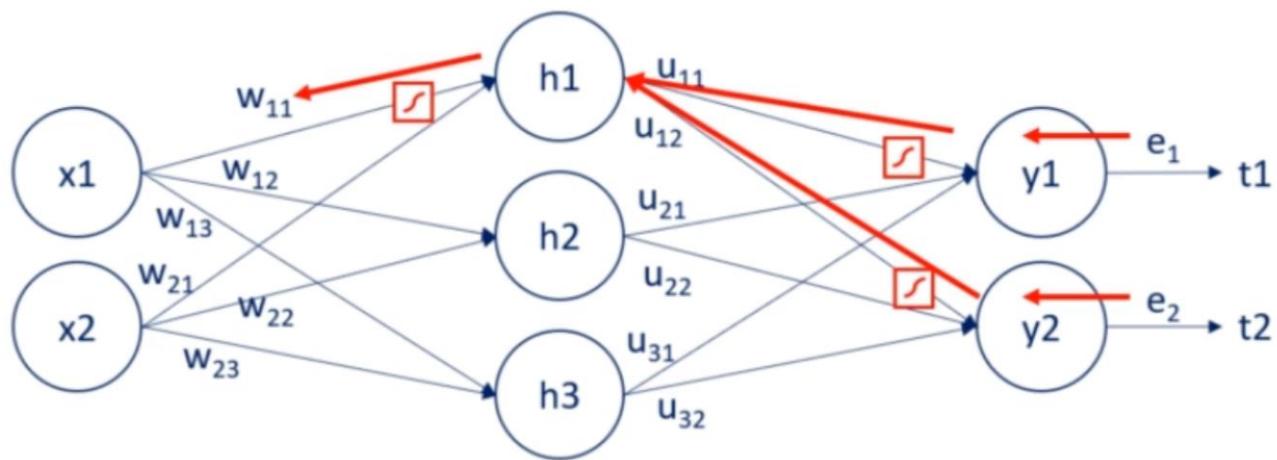
we find its derivative (gradient)  $\mathbf{u}_{i+1} = \mathbf{u}_i - \eta \nabla_{\mathbf{u}} L(\mathbf{u}_i)$

and update the coefficient.

# Backpropagation

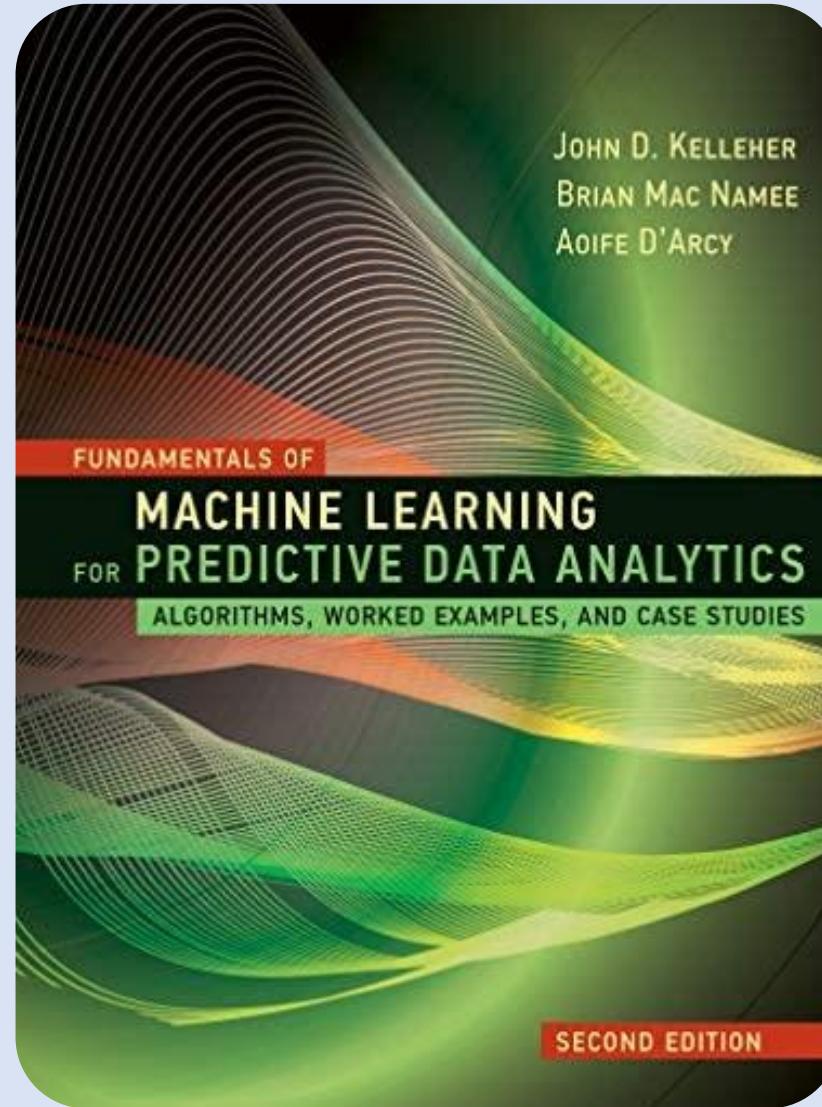


# Backpropagation mathematics



$L^2$  –Norm Loss and Sigmoid Activation

$$\frac{\partial L}{\partial w_{ij}} = \delta_j x_i, \text{ where } \delta_j = \sum_k \delta_k w_{jk} y_j (1 - y_j)$$



## CONTENTS

Series Foreword vii

Preface ix

Acknowledgments xi

- 1 Introduction to Deep Learning 1
- 2 Conceptual Foundations 39
- 3 Neural Networks: The Building Blocks of Deep Learning 65
- 4 A Brief History of Deep Learning 101
- 5 Convolutional and Recurrent Neural Networks 159
- 6 Learning Functions 185
- 7 The Future of Deep Learning 231

Glossary 251

Notes 257

References 261

Further Readings 267

Index 269

# DEEP LEARNING

JOHN D. KELLEHER



THE MIT PRESS ESSENTIAL KNOWLEDGE SERIES