# Natural Language Processing

**Lecture 15**

Dirk Hovy

dirk.hovy@unibocconi.it
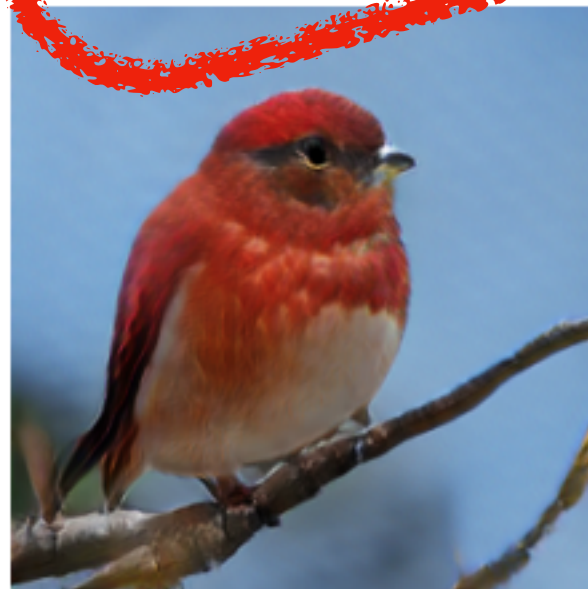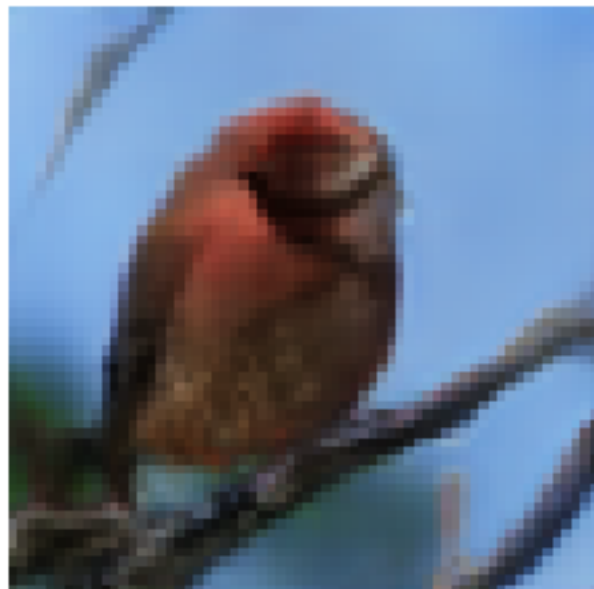
@dirk_hovy

Bocconi

# Neural Nets Everywhere
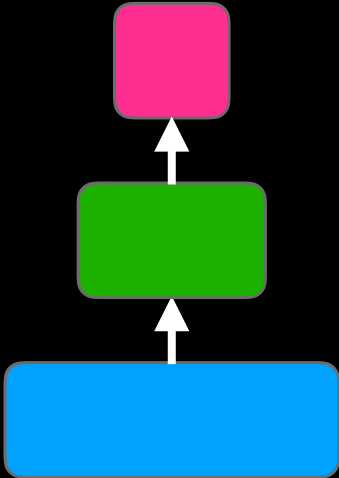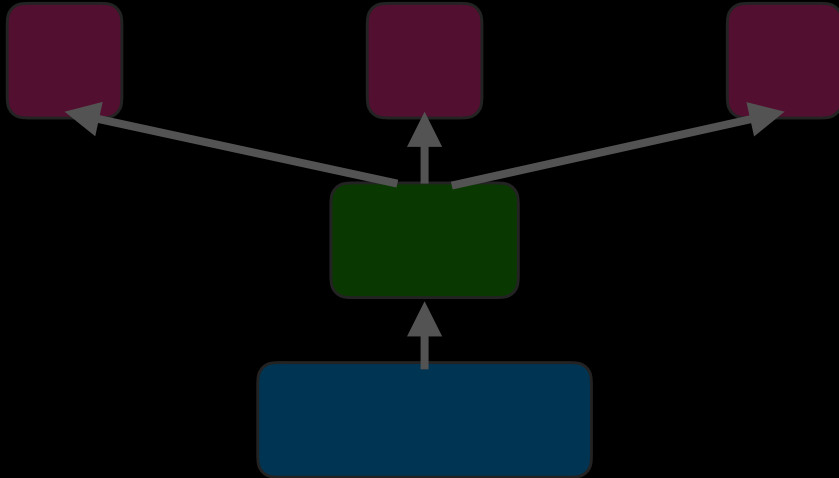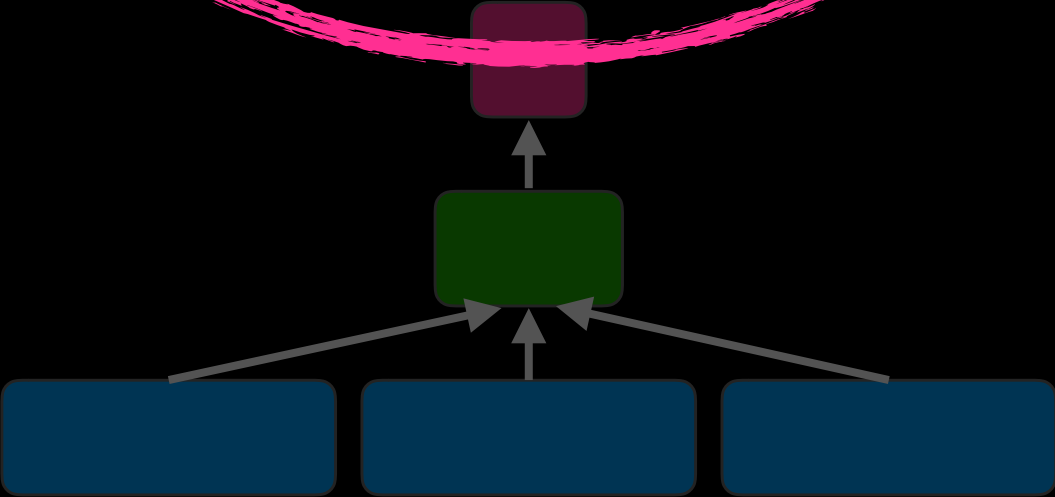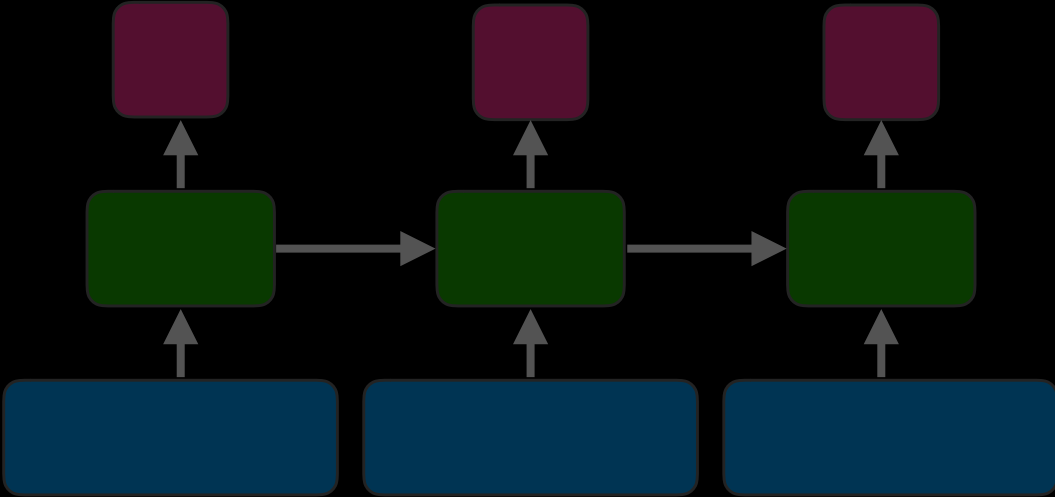
this bird is red with white and has a very short beak

FAKE NEWS

Bocconi

# Goals for Today

- Learn the basic difference between **neural architectures**

- Understand the **perceptron** as a basic element

- Understand training through **backpropagation**
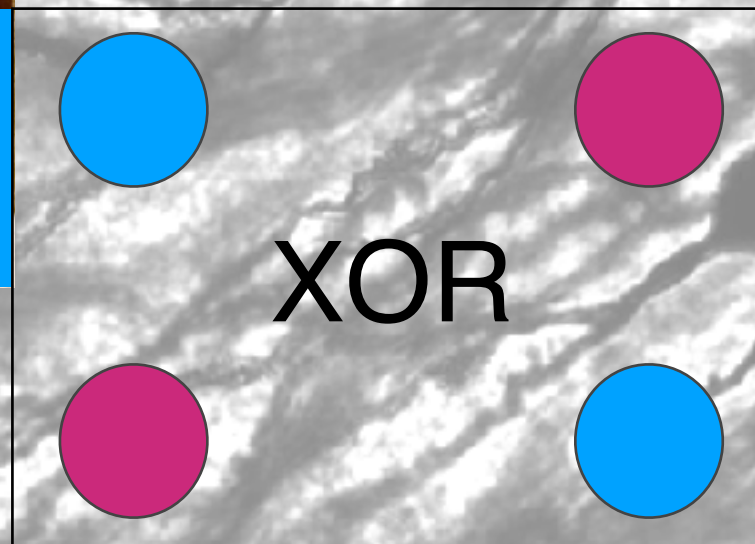
- Learn about **dropout regularization**

**Bocconi**

# Types of Neural Models

|  | Fixed length | Variable length |
|---|---|---|
| **Fixed length** | Logistic Regression, Perceptron, Feed-Forward Network, Deep Belief Network… | Multitask Learning, Decoder |
| **Variable length** | Convolutional Neural Networks (CNN) | Recurrent Neural Networks (RNN), Hidden Markov Models (HMM), Conditional Random Fields |

The Perceptron can learn anything!!!
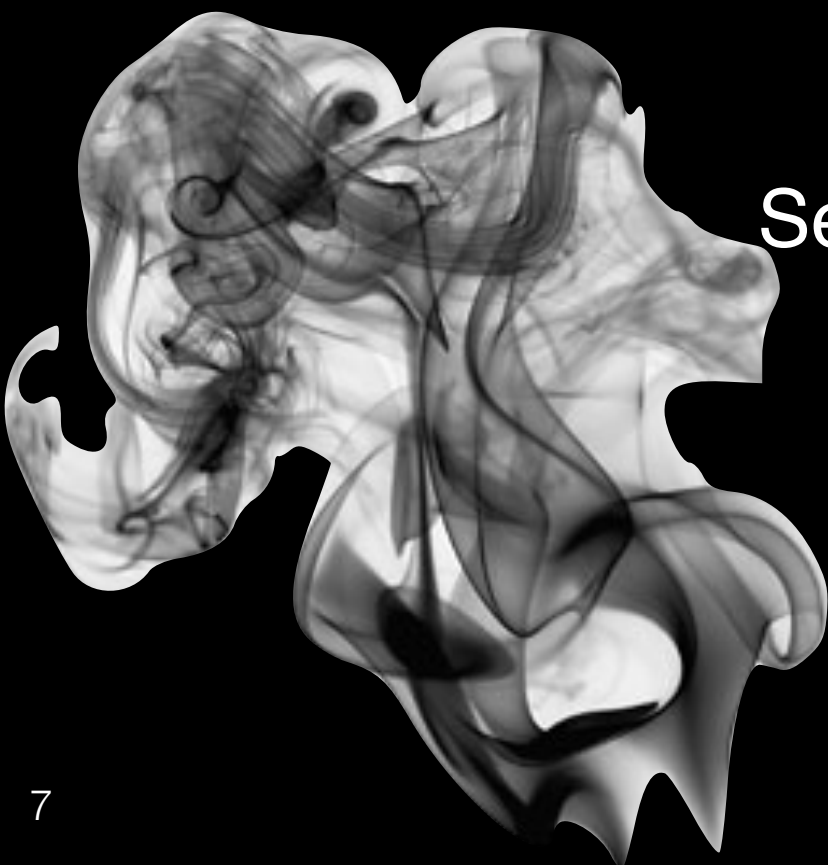
Frank Rosenblatt (1928–1971)

XOR

The Perceptron fails at learning even basic concepts

Marvin Minsky (1927–2016)

# The Perceptron

# A Threshold Unit

Sensor array

$$\Sigma$$

if > threshold

Total smoke

Bocconi

# OR-Perceptron

**INPUT**

$x_1$

1

**ACTIVATION**

**WEIGHTS**

$y$

1

$x_2$

$$f(X) = w_1 x_1 + w_2 x_2$$

$$\hat{y} = \begin{cases} +1 & if \ f(X) \geq \ 1 \\ -1 & otherwise \end{cases}$$

| x₁ | x₂ | y |
|----|----|----|
| 0 | 0 | -1 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 1 |

Bocconi

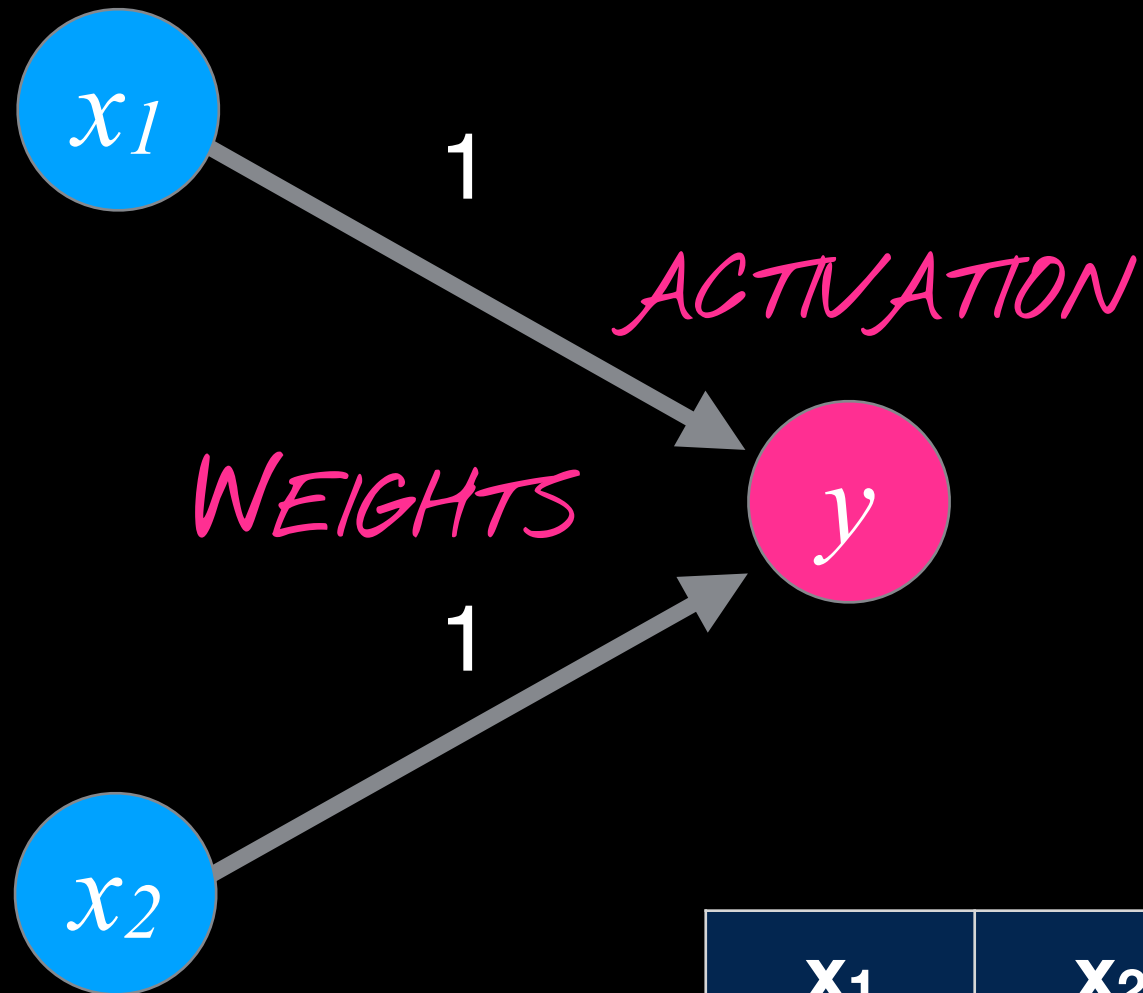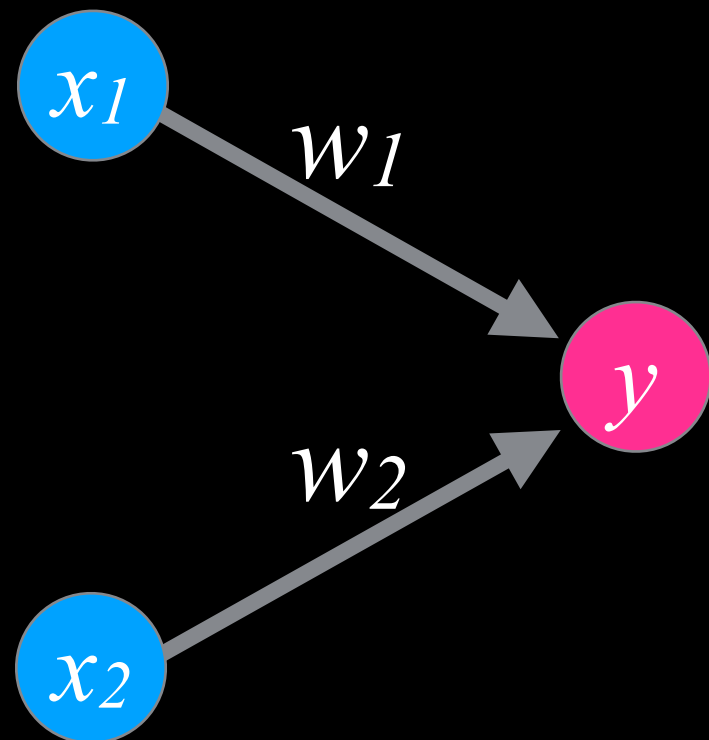# AND-Perceptron

INPUT

$x_1$

1

ACTIVATION

WEIGHTS

$y$

1

$x_2$

$$f(X) = w_1 x_1 + w_2 x_2$$

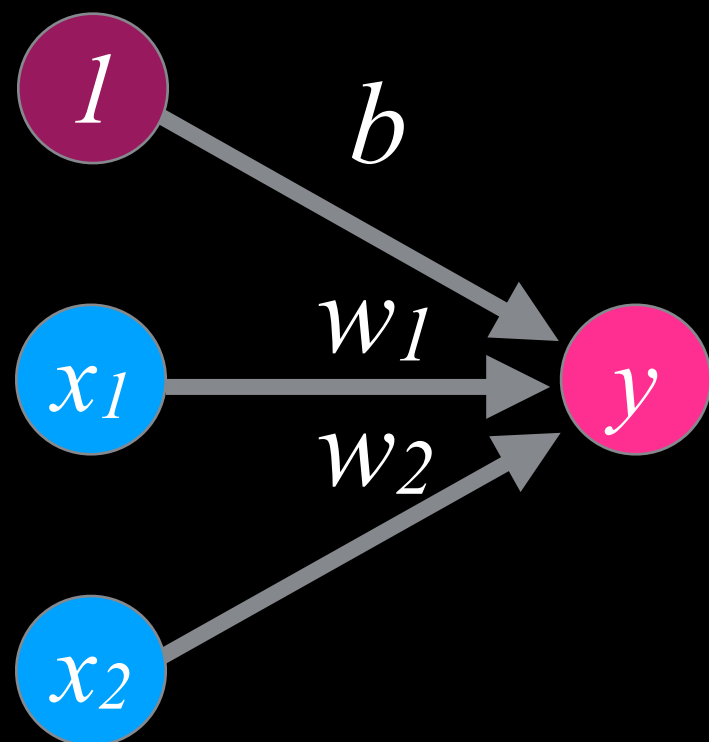$$\hat{y} = \begin{cases} +1 & if \ f(X) \geq 2 \\ -1 & otherwise \end{cases}$$

| $x_1$ | $x_2$ | y |
|---|---|---|
| 0 | 0 | -1 |
| 1 | 0 | -1 |
| 0 | 1 | -1 |
| 1 | 1 | 1 |

Bocconi

# Learn the Threshold

$$f(X) = w_1 x_1 + w_2 x_2$$

$$\hat{y} = \begin{cases} +1 & if \ f(X) \geq t \\ -1 & otherwise \end{cases}$$

$$f(X) = w_1 x_1 + w_2 x_2 + \textbf{\textit{b}}$$

$$\hat{y} = \begin{cases} +1 & if \ f(X) \geq \textbf{\textit{0.5}} \\ -1 & otherwise \end{cases}$$

Bocconi

# Learning to Distinguish

$f(X) = 1\,x_1 + 1\,x_2 - 1$

$f(X) = 1\,x_1 + 1\,x_2 - 2$
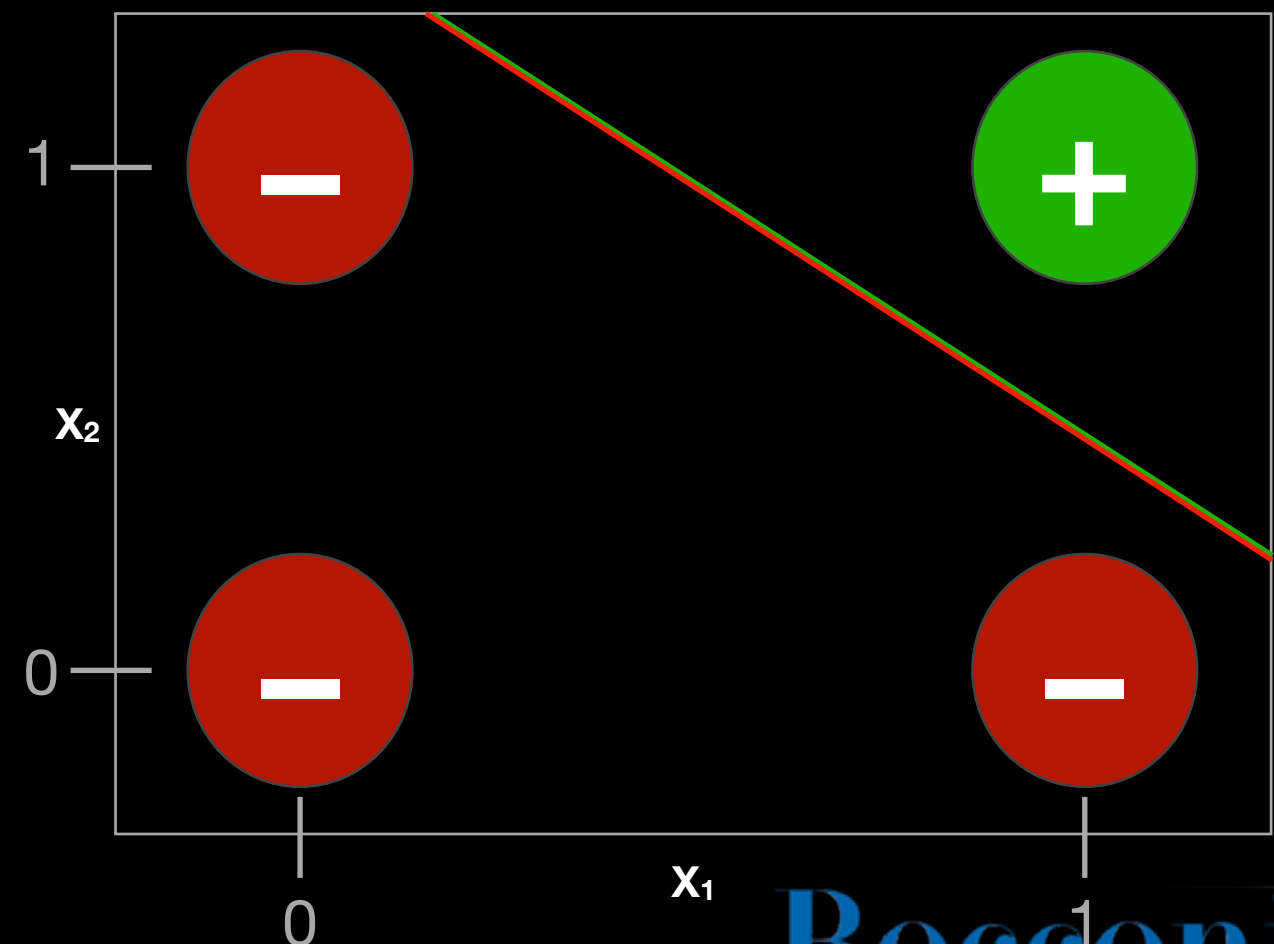
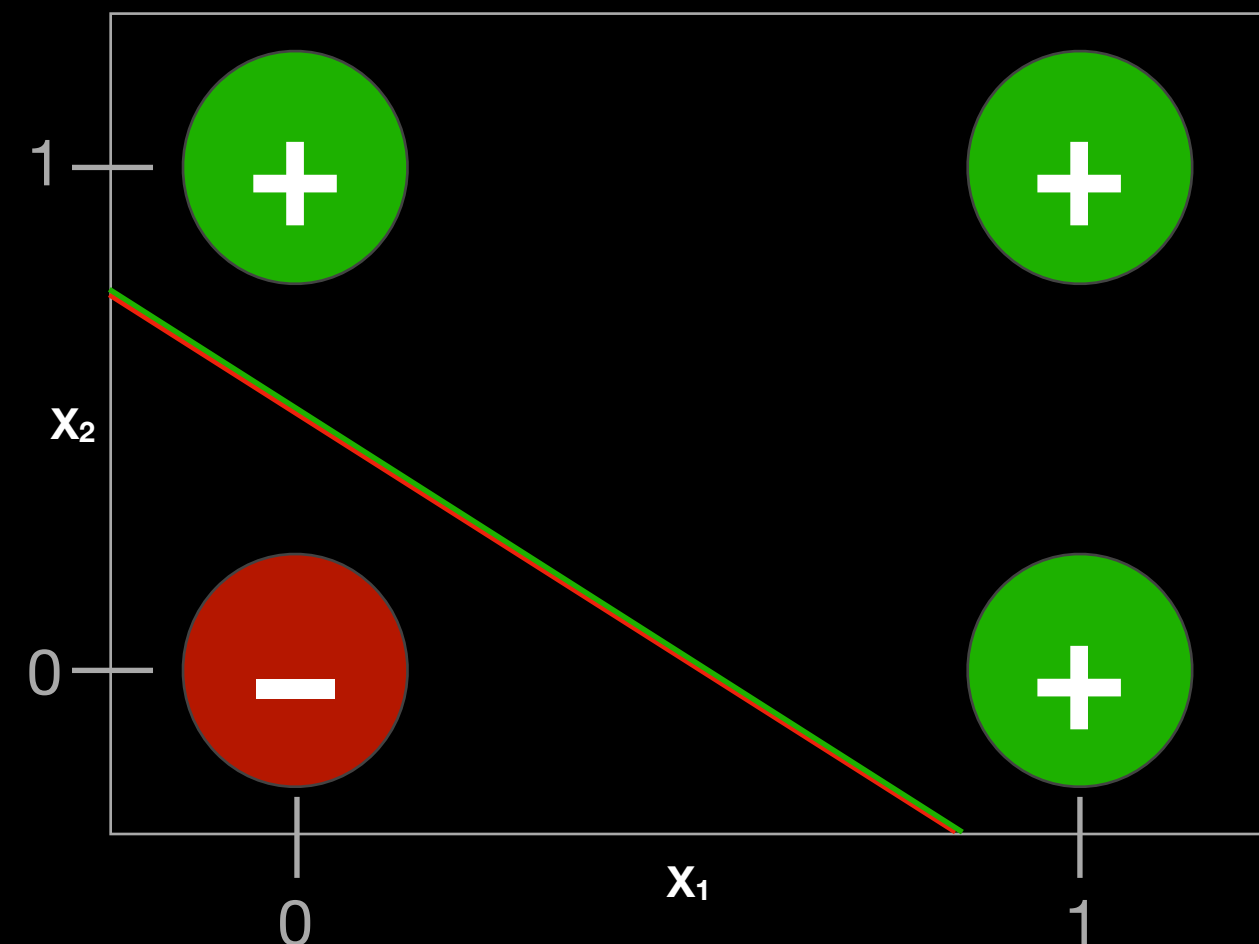OR

| $x_1$ | $x_2$ | y |
|-------|-------|-----|
| 0 | 0 | -1 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 1 |

AND

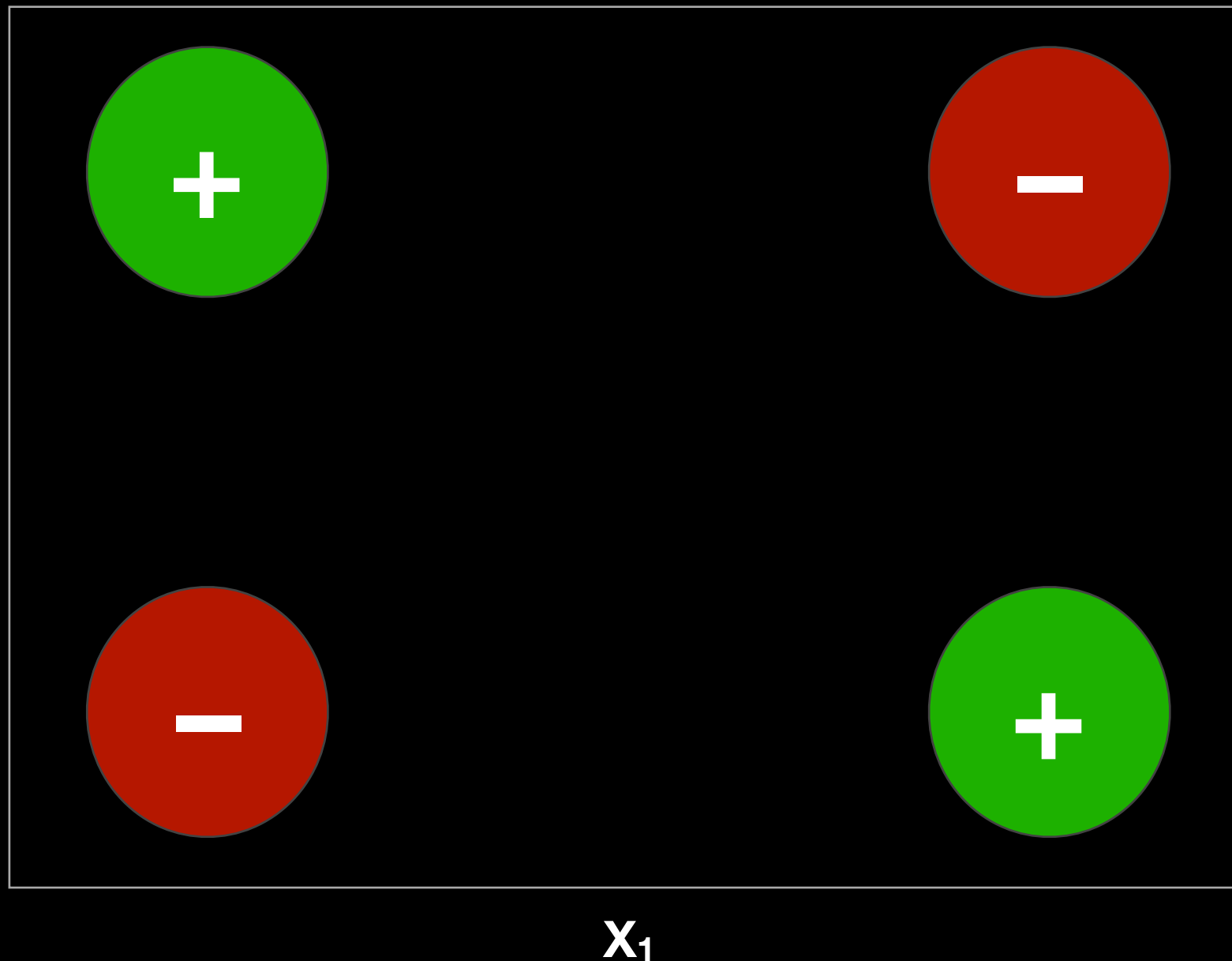| $x_1$ | $x_2$ | y |
|-------|-------|-----|
| 0 | 0 | -1 |
| 1 | 0 | -1 |
| 0 | 1 | -1 |
| 1 | 1 | 1 |

# The XOR Limit

| x₁ | x₂ | y |
|----|----|----|
| 0 | 0 | -1 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | -1 |

*Linearize this!*

Marvin Minsky
(1927–2016)

# Step 1: Non-Linearity

Bocconi

# Nonlinear Activation Functions
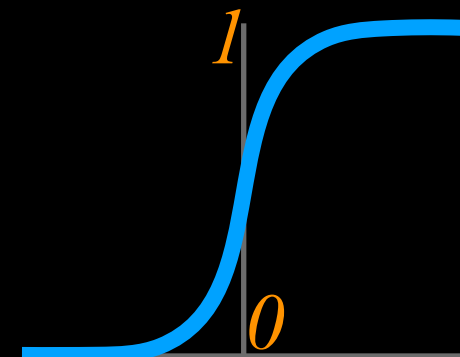
$$f(X) = \boldsymbol{a(}w_1 x_1 + w_2 x_2 + b\boldsymbol{)}$$



SIGMOID
(S-LIKE)

**Logistic**
$$s(x) = \frac{1}{1+e^{-x}}$$

**tanh**
$$tanh(x)$$

**ReLU**
$$max(0, x)$$

Bocconi

# Step 2:
# Going Deep –
# The Multilayer Perceptron

Bocconi

# Multilayer Perceptron

$$f(X) = a_2(W^2 \, a_1(W^1 X + b^1) + b^2))$$



1

1

$b^2_1$

$b_{1,1}$

$b_{1,3}$   $b_{1,2}$

$h_1$

$w^2_1$

$w_{1,1}$

$x_1$

$w_{1,2}$

$w_{1,3}$

$h_2$

$w^2_2$

$y$

$w_{2,1}$

$w_{2,2}$

$x_2$

$w_{2,3}$

$h_3$

$w^2_3$

Hidden

Layer

- How many perceptrons do you see?

- How many parameters?

Can Approximate

Any Function!

Bocconi

# The XOR Limit

# Multi-Class Output



$$f_i(X) = a(w_{1,i} x_1 + w_{2,i} x_2 + b_i)$$

$b_1$

$1$

$b_3$  $b_2$

$y_1$  positive

$w_{1,1}$

$x_1$  $w_{1,2}$  $y_2$  negative  $\hat{y} = \underset{i}{argmax}\, f_i(X)$

$w_{1,3}$

$w_{2,1}$

$w_{2,2}$  $y_3$  neutral

$x_2$  $w_{2,3}$

# Enter the Matrix



$$y_i = a(w_{1,i}x_1 + w_{2,i}x_2 + b_i)$$

$$Y = a(\ W^i \quad X + b^i)$$

# Learning

# Decision Boundary

book (+1)

magazine (–1)

thickness

size

# Error!



**neutral** **positive**

$-\sum log(\hat{y}) \bullet y$

*How far off are we?*

$1$   $b_1$   $b_2$   $b_3$

$y_1$   $1.0$   $-1.2$

$x_1$   $w_{1,1}$   $w_{1,2}$   $w_{1,3}$

$y_2$   $0.0$   $+$   $0.0$   $=$   $1.2$

$x_2$   $w_{2,1}$   $w_{2,2}$   $w_{2,3}$

$y_3$   $0.0$   $+$   $0.0$

*Cross-Entropy*

Bocconi

# Backpropagation



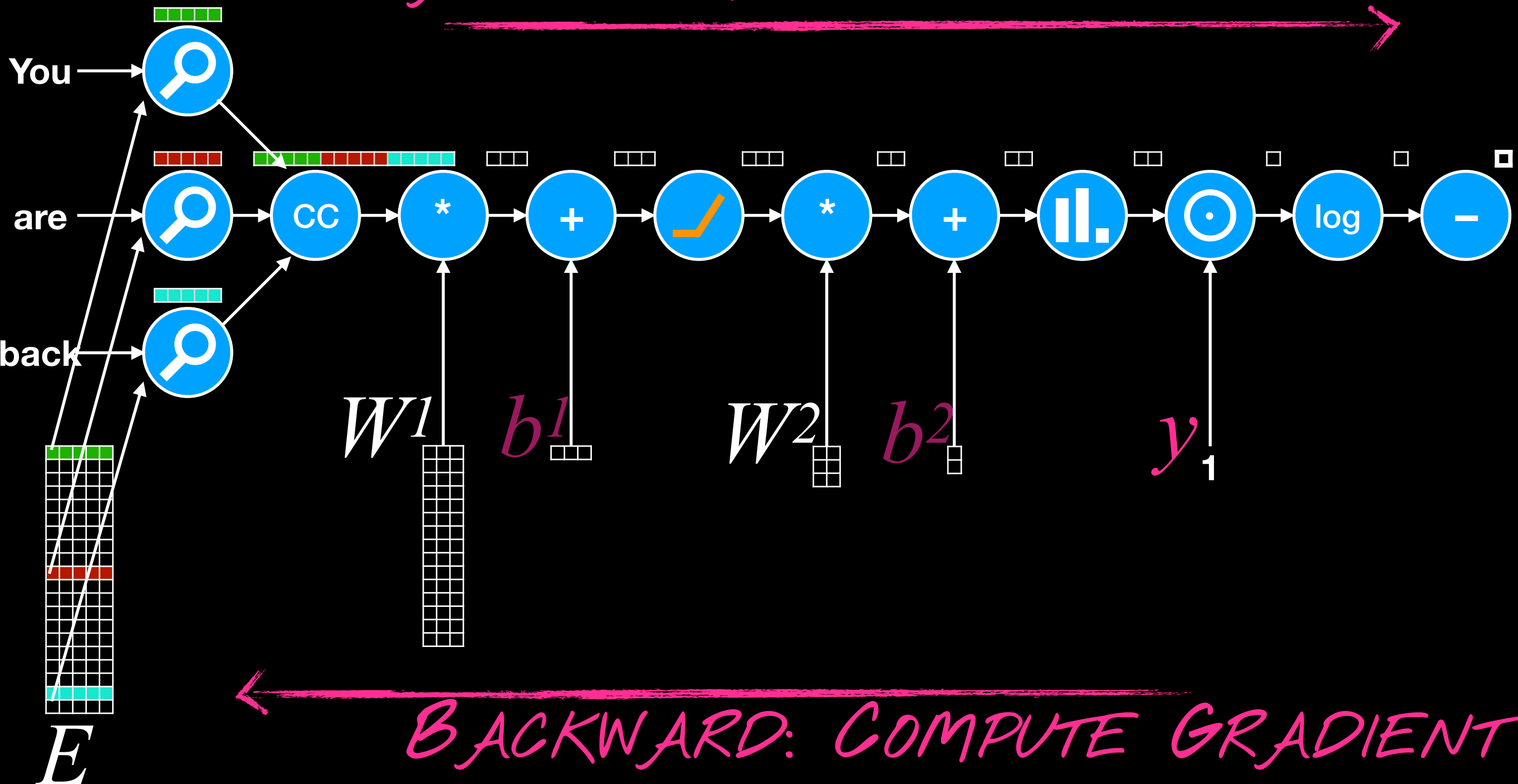- Adjust weights/bias proportionately to change, using **Stochastic Gradient Descent**

- In deeper networks: compute effect on previous layer activation, then adjust *their* incoming weights accordingly, using **Chain Rule**

- If input layer is adjusted as well = learning representations

# Computational Graph



FORWARD: COMPUTE ERROR

You

are

back

cc * + ⟋ * + ‖. ⊙ log −

$W^1$  $b^1$  $W^2$  $b^2$  $y_1$

$E$

BACKWARD: COMPUTE GRADIENT

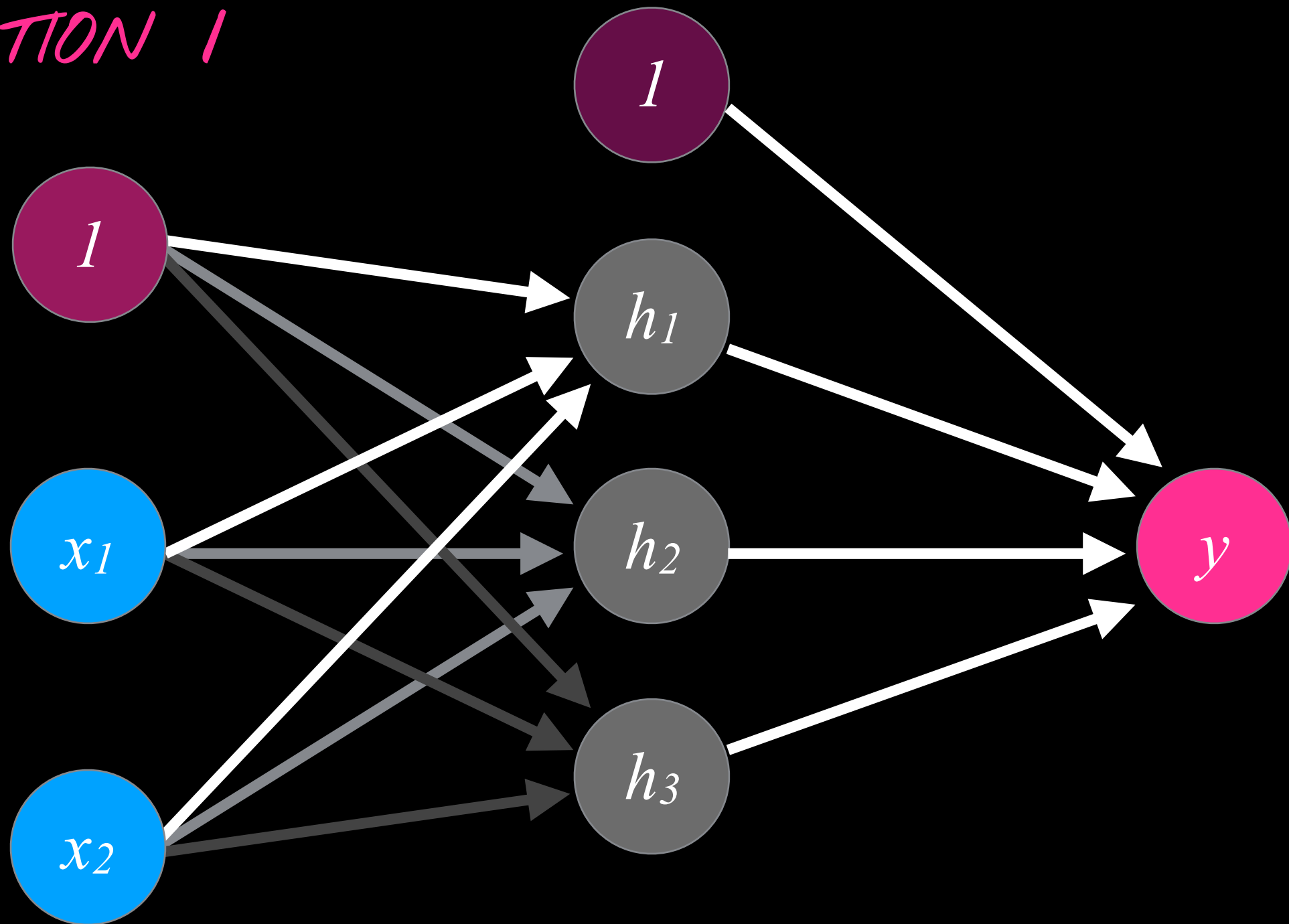Bocconi

# Regularization with Dropout

# Overfitting

DITCH

NO
DITCH

- ALVINN autonomous vehicle, trained on a particular stretch of road

- Drove off the road when turned around => focused on having a ditch on one side

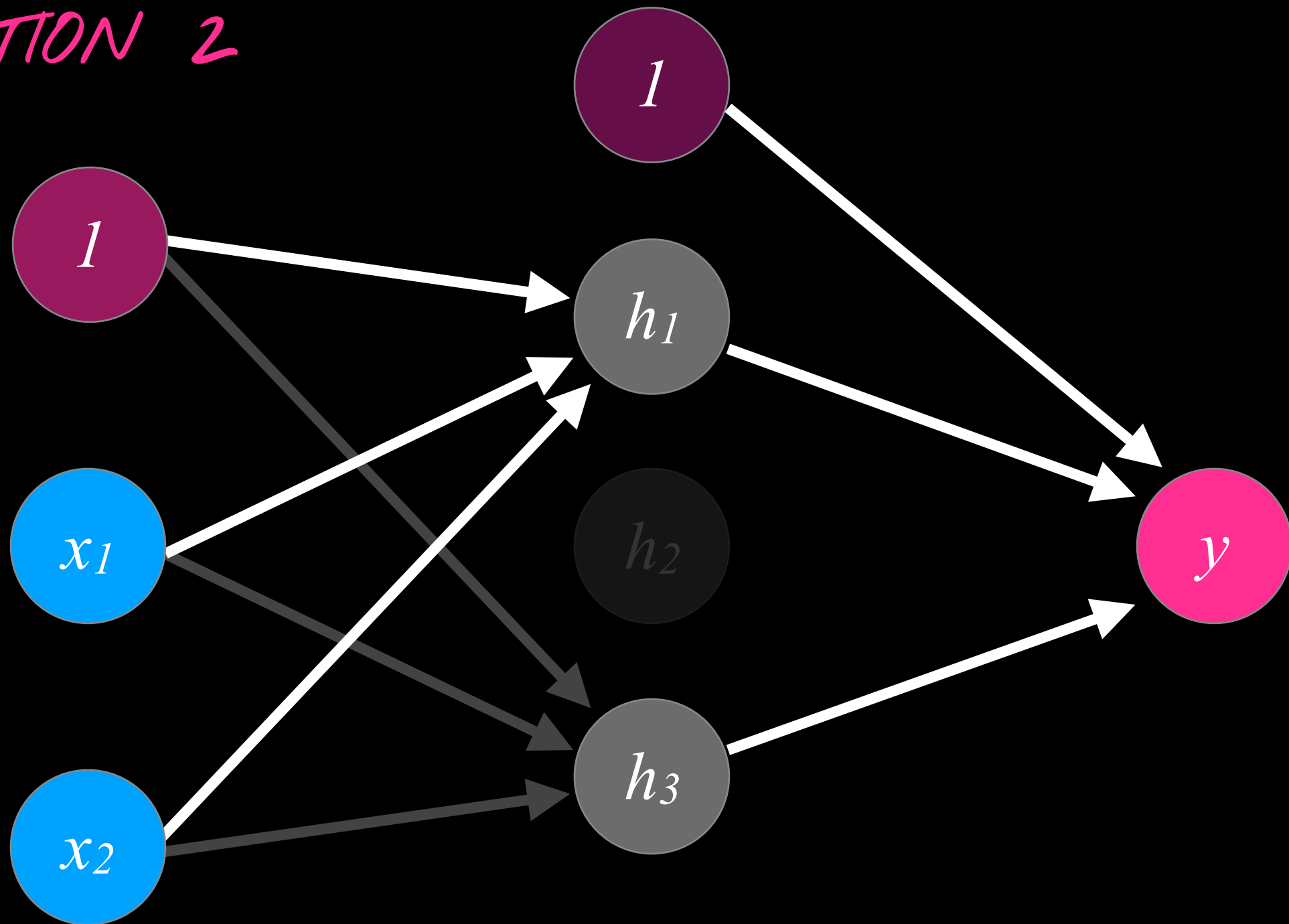- Idea: randomly remove nodes to avoid over-reliance
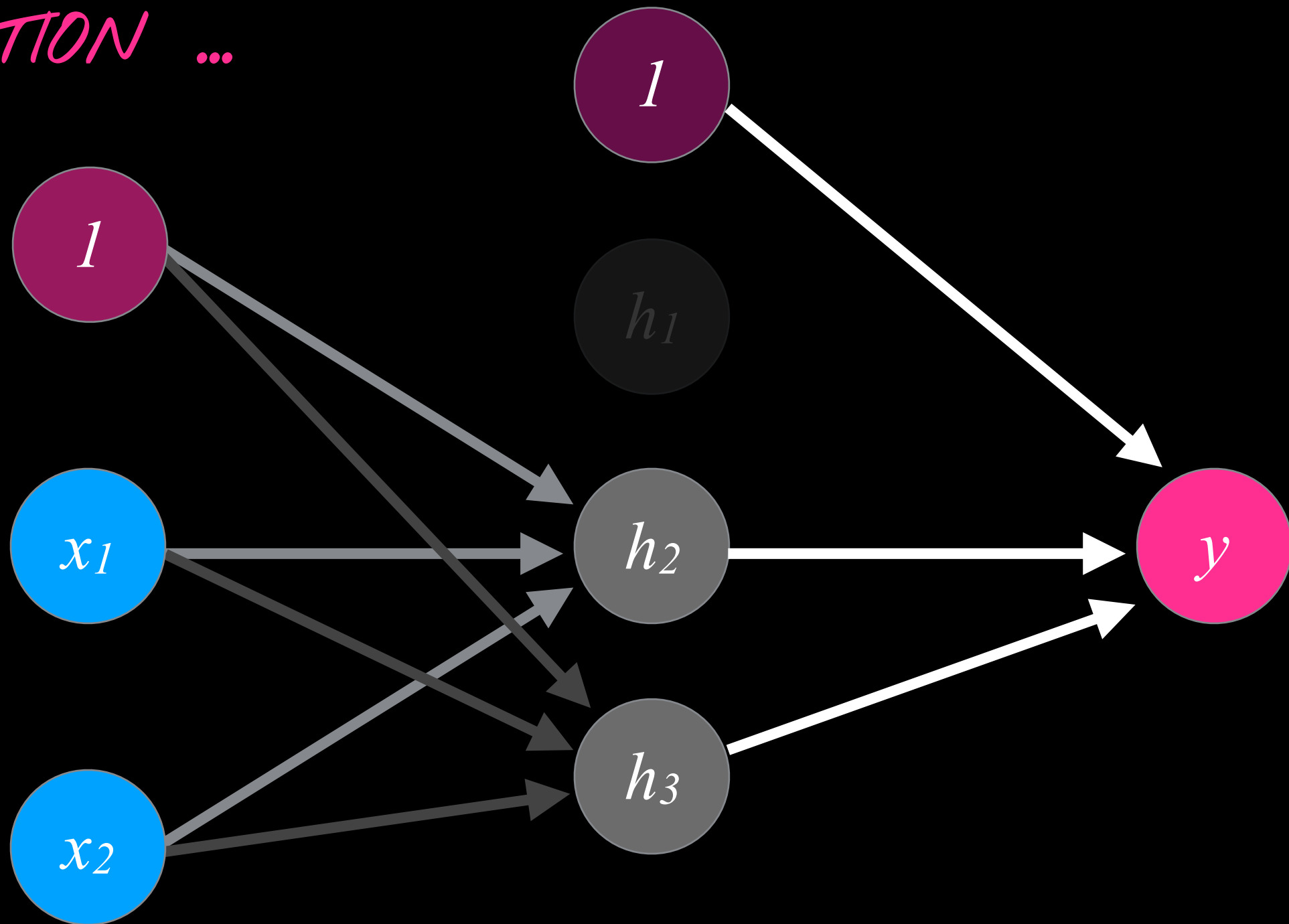
Bocconi

# Dropout



ITERATION 1

# Dropout

# Dropout



*Iteration ...*

# Wrapping up

# Take Home Points

- The **perceptron** is the basic building block of NNs

- Several perceptrons are a **Multilayer Perceptron** or **Feedforward Network**

- Each layer is matrix multiplication wrapped in an **activation function** (usually **ReLU**)

- Training **backpropagates** an error through the network to change weights

- **Dropout** helps regularize networks by randomly deleting nodes

*Bocconi*

# Moar Sources

- 3Blue1Brown: https://www.youtube.com/playlist?list=PLZHQObOWTQDNU6R1_67000Dx_ZCJB-3pi

- Yoav Goldberg Primer: https://arxiv.org/pdf/1510.00726.pdf

- The Keras book

- …

Bocconi