



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 11: Neural Networks I

Slides adapted from Chris Ketelsen, Jordan Boyd-Graber,
and Noah Smith

Administrivia

- Final project team formation
- Proposal ideas
- HW3 released

Learning Objectives

- Understanding feed-forward neural networks

Outline

- Revisiting logistic regression
- Feed-forward neural networks

Outline

- Revisiting logistic regression
- Feed-forward neural networks

Revisiting Logistic Regression

$$P(\underline{Y} = 0 \mid \underline{x}, \beta) = \frac{1}{1 + \exp \left[\beta_0 + \sum_j \beta_j x_j \right]}$$

$$P(Y = 1 \mid \underline{x}, \beta) = \frac{\exp \left[\beta_0 + \sum_j \beta_j x_j \right]}{1 + \exp \left[\beta_0 + \sum_j \beta_j x_j \right]}$$

$$\mathcal{L} = - \sum_i \log P(y^{(i)} \mid \mathbf{x}^{(i)}, \beta)$$

Next slide $\beta_1 = [0, \dots, 0]$
 $\beta_2 = [\beta_0, \dots, \beta_d]$

$x \xrightarrow{\text{linear}} \begin{bmatrix} 0 \\ \beta_0 + \sum_j \beta_j x_j \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} P(y=0|x) \\ P(y=1|x) \end{bmatrix}$

$\exp(0)$

$\exp(0) + \exp(\beta_0 + \sum_j \beta_j x_j)$

$0 = \beta^{(0)} \cdot x$

$\beta \cdot x$

$\begin{pmatrix} \beta^{(0)} \\ \beta \end{pmatrix} x = \begin{bmatrix} 0 \\ \beta_0 + \sum_j \beta_j x_j \end{bmatrix}$

softmax $\begin{bmatrix} \frac{\exp 0}{\exp 0 + \exp(\beta_0 + \sum_j \beta_j x_j)} \\ \frac{\exp(\beta_0 + \sum_j \beta_j x_j)}{\exp 0 + \exp(\beta_0 + \sum_j \beta_j x_j)} \end{bmatrix}$

Revisiting Logistic Regression

- Transformation on x (we map class labels from $\{0, 1\}$ to $\{1, 2\}$):

$$\underline{l_k} = \beta_k^T \mathbf{x}, k = 1, 2$$

$$\underline{o_k} = \frac{\exp l_k}{\sum_{c \in \{1, 2\}} \exp l_c}, k = 1, 2$$

Revisiting Logistic Regression

- Transformation on \mathbf{x} (we map class labels from $\{0, 1\}$ to $\{1, 2\}$):

$$l_k = \beta_k^T \mathbf{x}, k = 1, 2 \quad \text{linear layer}$$

$$o_k = \frac{\exp l_k}{\sum_{c \in \{1, 2\}} \exp l_c}, k = 1, 2 \quad \text{softmax layer}$$

Revisiting Logistic Regression

- Transformation on \mathbf{x} (we map class labels from $\{0, 1\}$ to $\{1, 2\}$):

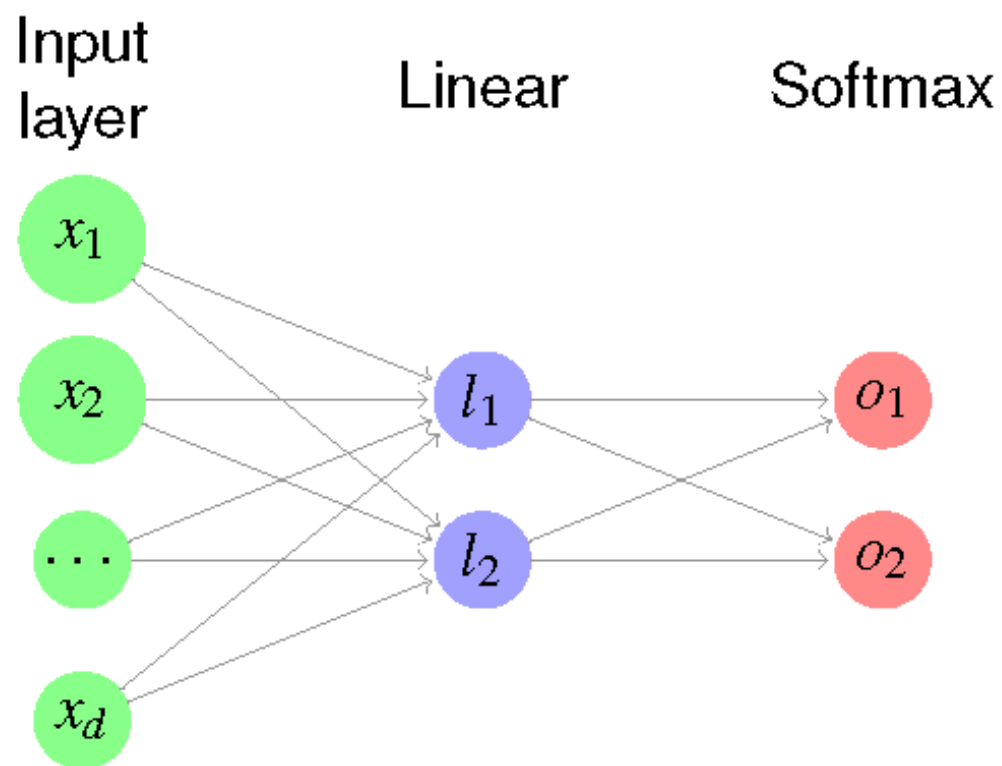
$$l_k = \beta_k^T \mathbf{x}, k = 1, 2 \quad \text{linear layer}$$

$$o_k = \frac{\exp l_k}{\sum_{c \in \{1, 2\}} \exp l_c}, k = 1, 2 \quad \text{softmax layer}$$

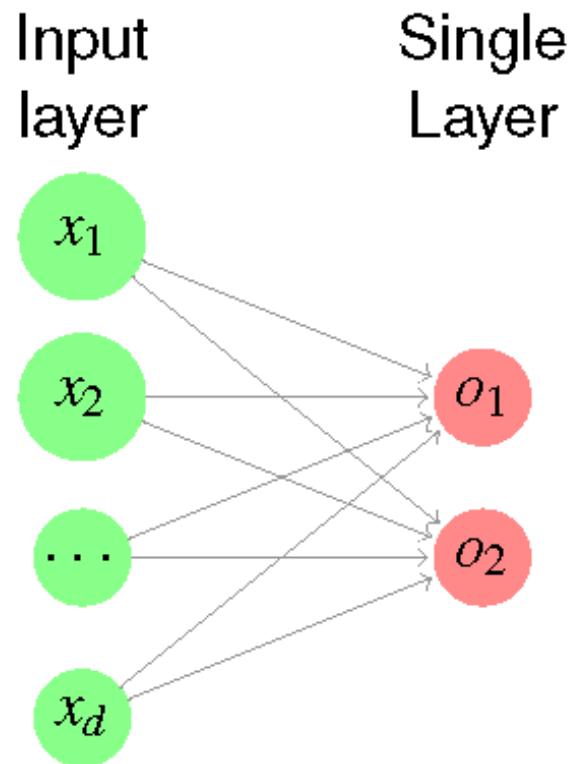
- Objective function (using cross entropy $-\sum_i p_i \log q_i$):

$$\mathcal{L}(Y, \hat{Y}) = - \sum_i \left[P(y^{(i)} = 1) \log P(\hat{y}_i = 1 \mid \mathbf{x}^{(i)}, \beta) + P(y^{(i)} = 0) \log \hat{P}(y_i = 0 \mid \mathbf{x}^{(i)}, \beta) \right]$$

Logistic Regression as a Single-layer Neural Network



Logistic Regression as a Single-layer Neural Network



$$\frac{\exp(x)}{\sum_v \exp(x_v)}$$

k

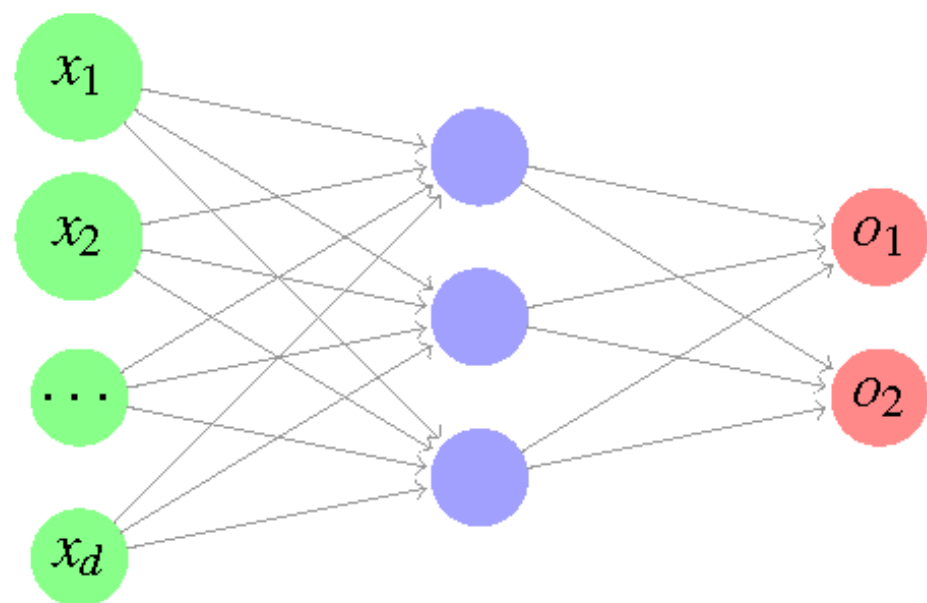
Outline

- Revisiting logistic regression
- **Feed-forward neural networks**

Deep Neural networks

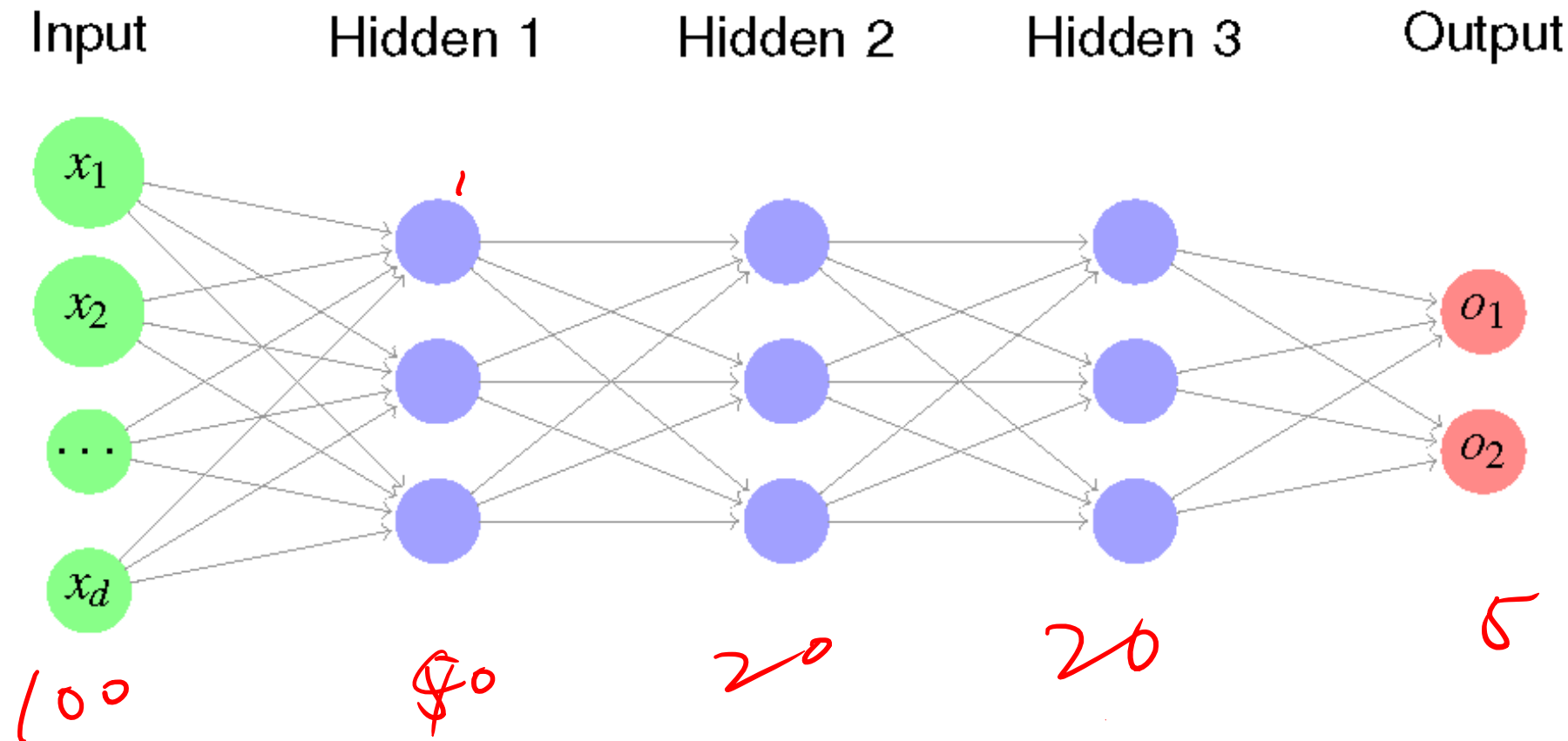
A two-layer example (one hidden layer)

Input Hidden Output



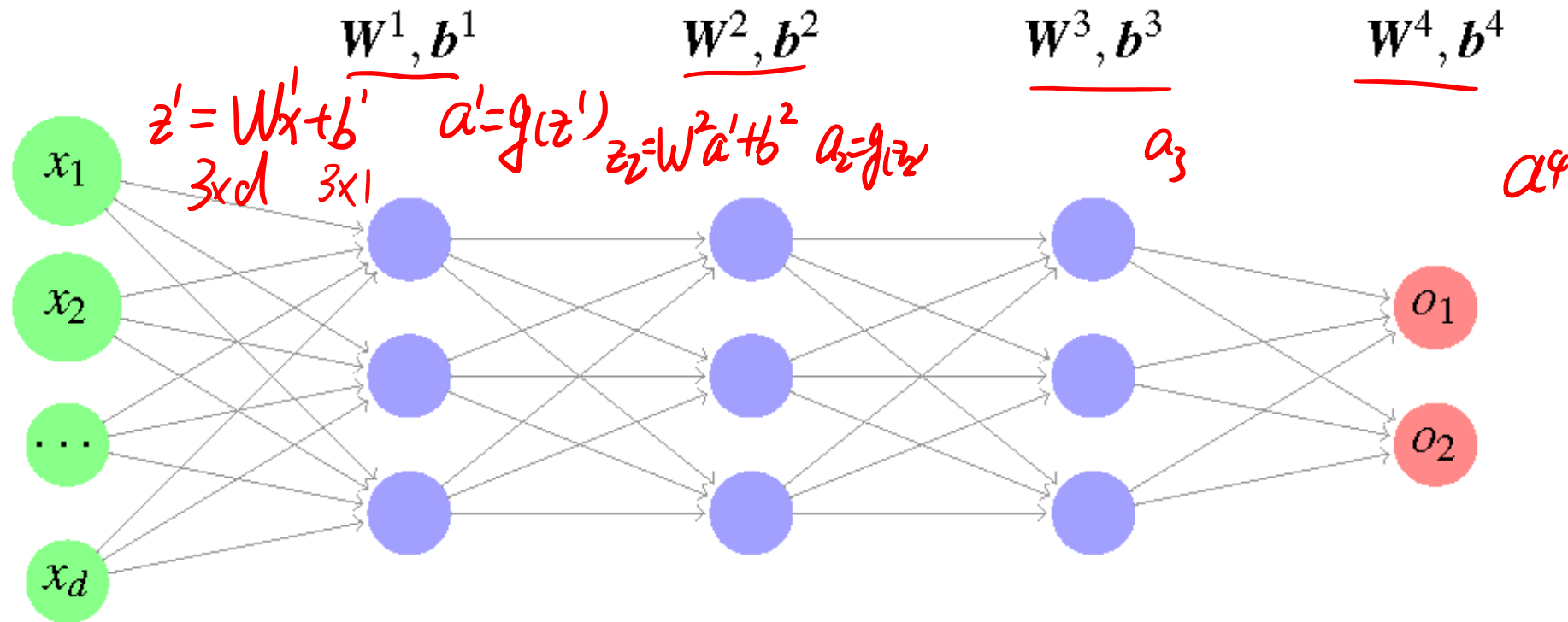
Deep Neural networks

More layers:



Forward propagation algorithm

How do we make predictions based on a multi-layer neural network?
Store the biases for layer l in b^l , weight matrix in W^l



Forward propagation algorithm

Suppose your network has L layers
Make prediction for an instance x

- 1: Initialize $a^0 = x$
- 2: **for** $l = 1$ to L **do**
- 3: $z^l = W^l a^{l-1} + b^l$
- 4: $a^l = g(z^l)$
- 5: **end for**
- 6: The prediction \hat{y} is simply a^L

Nonlinearity

What happens if there is no nonlinearity?

$$\begin{aligned} & W^2(W'x + b') + b^2 \\ &= \underbrace{W^2 W'}_w x + \underbrace{W^2 b'}_b + b^2 \end{aligned}$$

Nonlinearity

What happens if there is no nonlinearity?

Linear combinations of linear combinations are still linear combinations.

Neural networks in a nutshell

- Training data $S_{\text{train}} = \{(\mathbf{x}, y)\}$
- Network architecture (model)
- Loss function (objective function)
- Learning (next week)

$$\hat{y} = f_w(\mathbf{x})$$

g

$$\mathcal{L}(y, \hat{y})$$

*\downarrow
 a^L*

Nonlinearity Options

- Sigmoid

$$f(x) = \frac{1}{1 + \exp(x)}$$

- tanh

$$f(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

- ReLU (rectified linear unit)

$$f(x) = \max(0, x)$$

- softmax

$$x = \frac{\exp(\mathbf{x})}{\sum_{x_i} \exp(x_i)}$$

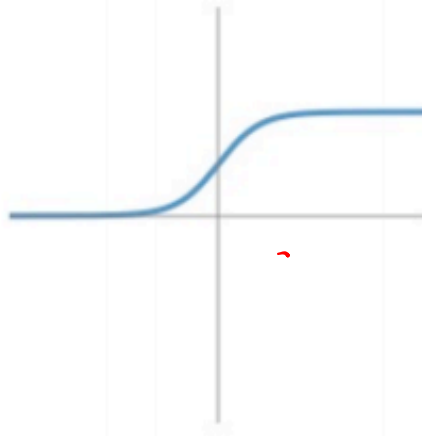
<https://keras.io/activations/>

Nonlinearity Options

$I(a > 0)$



Sigmoid



tanh

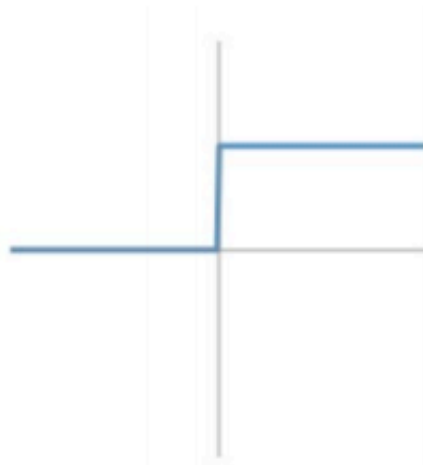


ReLU
 $\max(0, x)$

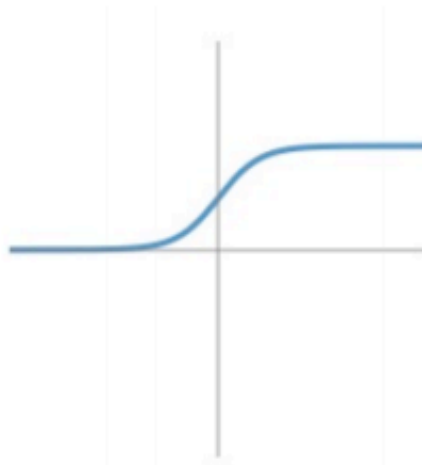


Nonlinearity Options

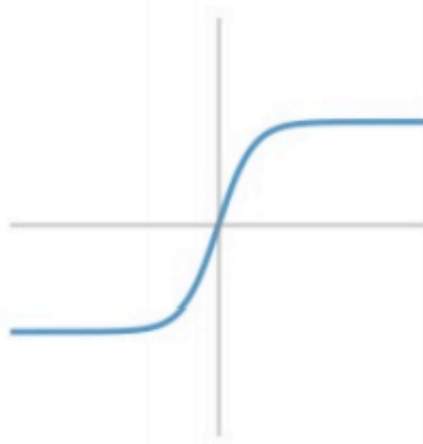
Perceptron



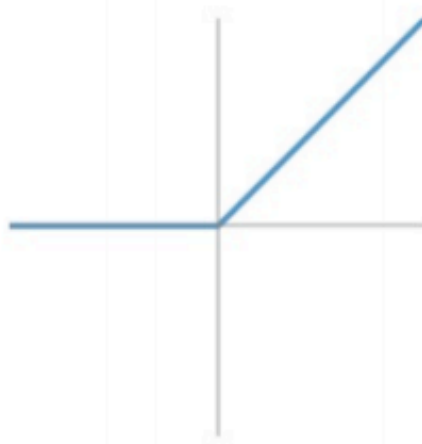
Sigmoid



Tanh



ReLU



Loss Function Options

- ℓ_2 loss

$$\sum_i (y_i - \hat{y}_i)^2$$

- ℓ_1 loss

$$\sum_i |y_i - \hat{y}_i|$$

- Cross entropy (logistic regression)

$$-\sum_i y_i \log \hat{y}_i$$

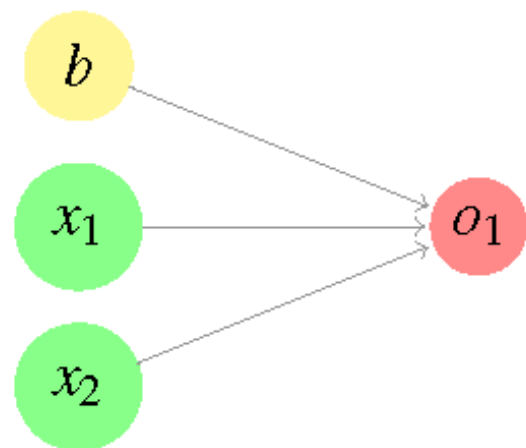
- Hinge loss (more on this during SVM)

$$\max(0, 1 - y\hat{y})$$

<https://keras.io/losses/>

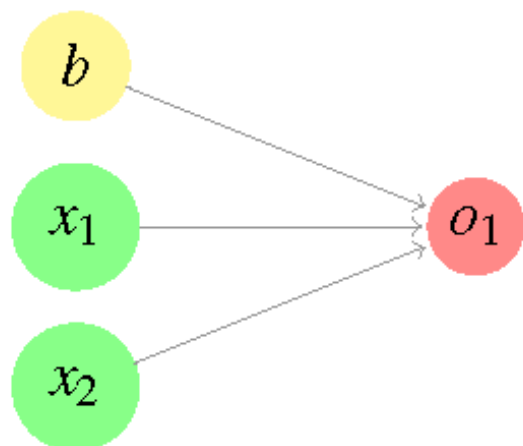
A Perceptron Example

$$\mathbf{x} = (x_1, x_2), y = f(x_1, x_2)$$



A Perceptron Example

$$\mathbf{x} = (x_1, x_2), y = f(x_1, x_2)$$



We consider a simple activation function

$$f(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

A Perceptron Example

Simple Example: Can we learn OR?

	<u>0</u>	<u>1</u>	<u>1</u>	<u>2</u>
x_1	0	1	0	1
x_2	0	0	1	1
$y = x_1 \vee x_2$	0	1	1	1

w_1, w_2 b 1

w_1, w_2
 $[1, 1]$ $-\frac{1}{2}$

$$w \quad b \quad \hat{y} = \begin{cases} 1 & w_1 x_1 + w_2 x_2 + b \geq 0 \\ 0 & w_1 x_1 + w_2 x_2 + b < 0 \end{cases}$$

if $(x_1 \text{ OR } x_2) = 1$

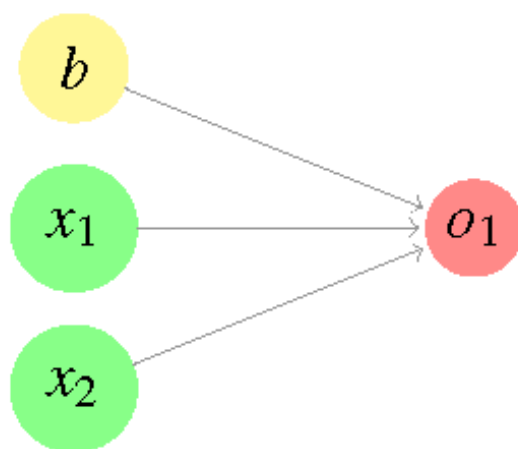
if $(x_1 \text{ OR } x_2) = 0$

A Perceptron Example

Simple Example: Can we learn OR?

x_1	0	1	0	1
x_2	0	0	1	1
$y = x_1 \vee x_2$	0	1	1	1

$$\mathbf{w} = (1, 1), b = -0.5$$



A Perceptron Example

Simple Example: Can we learn AND?

	<i>0</i>	<i>1</i>	<i>1</i>	<i>2</i>
x_1	0	1	0	1
x_2	0	0	1	1
$y = x_1 \wedge x_2$	0	0	0	1

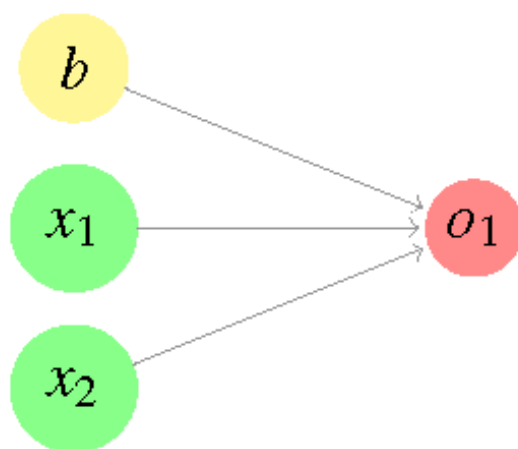
(1, 1) *-1.5*

A Perceptron Example

Simple Example: Can we learn AND?

x_1	0	1	0	1
x_2	0	0	1	1
$y = x_1 \wedge x_2$	0	0	0	1

$$\mathbf{w} = (1, 1), b = -1.5$$



A Perceptron Example

Simple Example: Can we learn NAND?

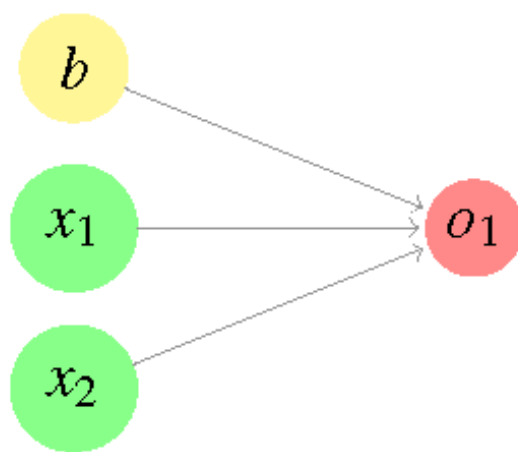
x_1	0	1	0	1
x_2	0	0	1	1
$y = \neg(x_1 \wedge x_2)$	1	1	1	0

A Perceptron Example

Simple Example: Can we learn NAND?

x_1	0	1	0	1
x_2	0	0	1	1
$y = \neg(x_1 \wedge x_2)$	1	1	1	0

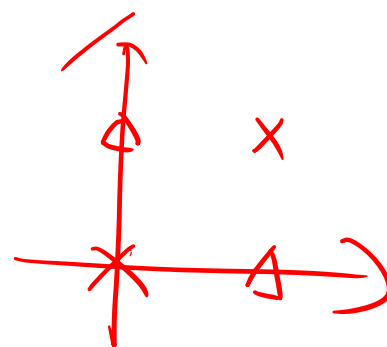
$$\mathbf{w} = (-1, -1), b = 1.5$$



A Perceptron Example

Simple Example: Can we learn XOR?

x_1	0	1	0	1
x_2	0	0	1	1
x_1 XOR x_2	0	1	1	0



A Perceptron Example

Simple Example: Can we learn XOR?

x_1		0	1	0	1	
x_2		0	0	1	1	
x_1	XOR	x_2	0	1	1	0

NOPE!

A Perceptron Example

Simple Example: Can we learn XOR?

x_1		0	1	0	1	
x_2		0	0	1	1	
x_1	XOR	x_2	0	1	1	0

NOPE!

But why?

A Perceptron Example

Simple Example: Can we learn XOR?

x_1		0	1	0	1	
x_2		0	0	1	1	
x_1	XOR	x_2	0	1	1	0

NOPE!

But why?

The single-layer perceptron is just a linear classifier, and can only learn things that are linearly separable.

A Perceptron Example

Simple Example: Can we learn XOR?

x_1		0	1	0	1	
x_2		0	0	1	1	
x_1	XOR	x_2	0	1	1	0

NOPE!

But why?

The single-layer perceptron is just a linear classifier, and can only learn things that are linearly separable.

How can we fix this?

A Perceptron Example

Increase the number of layers.

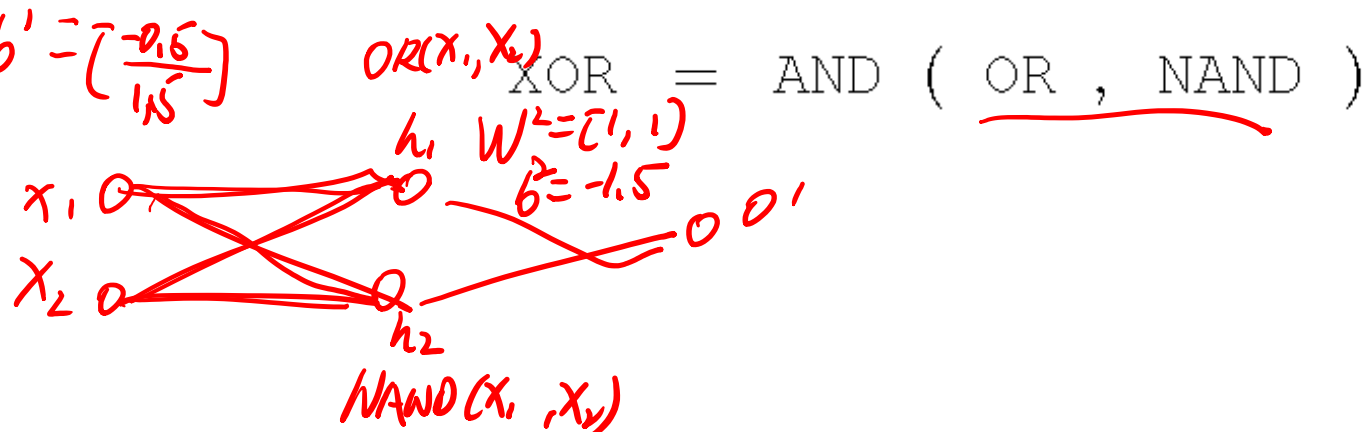
x_1		0	1	0	1	
x_2		0	0	1	1	
x_1	XOR	x_2	0	1	1	0

A Perceptron Example

Increase the number of layers.

2x2
 $W^1 = \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}$ OR
 NAND

$b^1 = \begin{bmatrix} -0.5 \\ 1.5 \end{bmatrix}$



x_1	0	1	0	1
x_2	0	0	1	1
x_1 XOR x_2	0	1	1	0

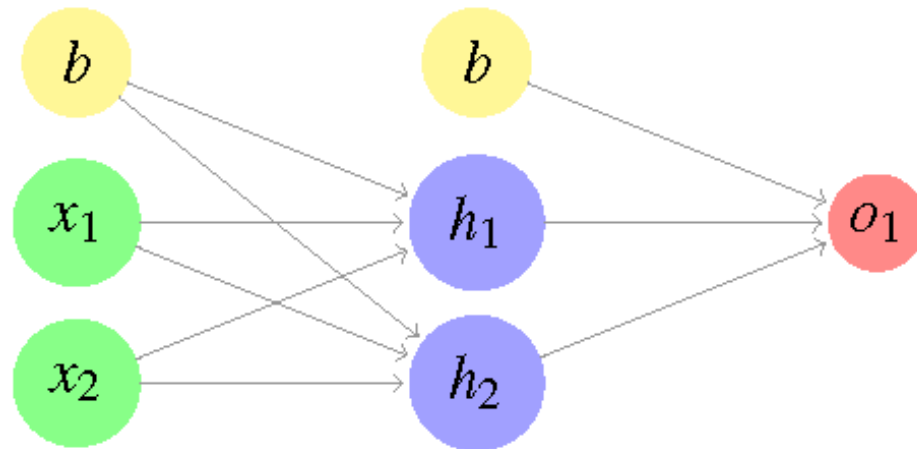
XOR = AND (OR , NAND)

A Perceptron Example

Increase the number of layers.

*Non-linear activations
are step functions.*

x_1	0	1	0	1
x_2	0	0	1	1
x_1 XOR x_2	0	1	1	0



$$W^1 = \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}, b^1 = \begin{bmatrix} -0.5 \\ 1.5 \end{bmatrix}$$

$$W^2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b^2 = -1.5$$

General Expressiveness of Neural Networks

Neural networks with a single hidden layer can approximate any measurable functions [Hornik et al., 1989, Cybenko, 1989].

Recap

- Logistic regression and perceptron can be seen as special cases of neural networks
- Feed-forward algorithm (forward propagation)