



Department of Computer Science  
UNIVERSITY OF COLORADO **BOULDER**



## Machine Learning: Chenhao Tan

University of Colorado Boulder  
LECTURE 14

Slides adapted from Jordan Boyd-Graber

## Logistics

---

- HW3 available on Github, due on October 16
- Prelim 1 grades have been released
- Project team match

## Outline

---

Revisiting Logistic Regression

Feed Forward Networks

## Revisiting Logistic Regression

---

$$P(Y = 0 \mid \mathbf{x}, \beta) = \frac{1}{1 + \exp \left[ \beta_0 + \sum_j \beta_j \mathbf{x}_j \right]}$$
$$P(Y = 1 \mid \mathbf{x}, \beta) = \frac{\exp \left[ \beta_0 + \sum_j \beta_j \mathbf{x}_j \right]}{1 + \exp \left[ \beta_0 + \sum_j \beta_j \mathbf{x}_j \right]}$$
$$\mathcal{L} = - \sum_i \log P(y^{(i)} \mid \mathbf{x}^{(i)}, \beta)$$

## Revisiting Logistic Regression

---

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

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

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

## 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  $\{1, 2\}$ ):

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

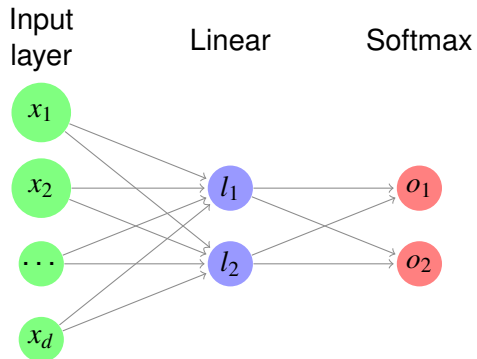
$$o_k = \frac{\exp l_k}{\sum_{c \in \{1, 2\}} \exp l_c}, k = 1, 2 \quad \text{softmax layer}$$

- 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

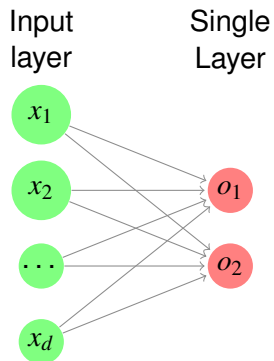
---





## Logistic Regression as a Single-layer Neural Network

---



## Outline

---

Revisiting Logistic Regression

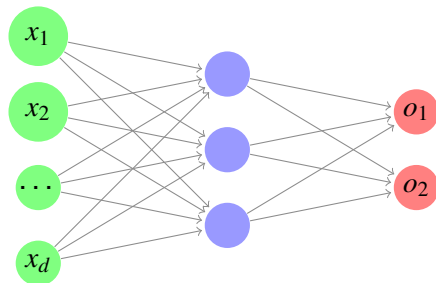
Feed Forward Networks

## Deep Neural networks

---

A two-layer example (one hidden layer)

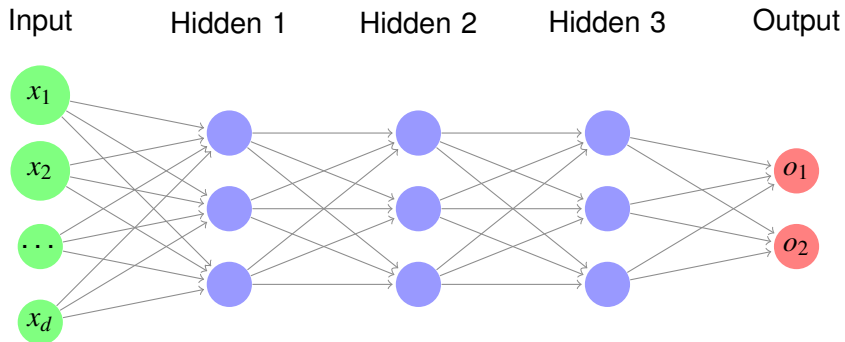
Input                  Hidden                  Output



## Deep Neural networks

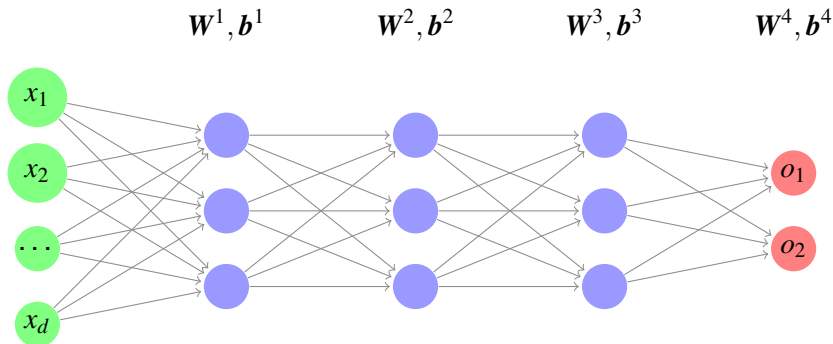
---

More layers:



## Forward propagation algorithm

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



## 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?

## 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)

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

- Loss function (objective function)

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

- Learning (next week)

## 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

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

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

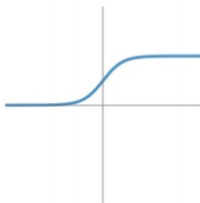
## Nonlinearity Options

---

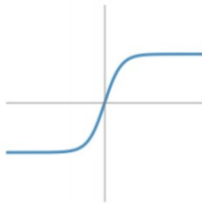
Perceptron



Sigmoid



Tanh



ReLU



## Loss Function Options

---

- $\ell_2$  loss

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

- $\ell_1$  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)

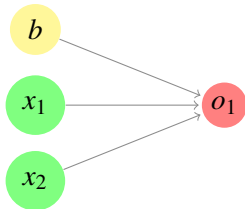
$$\max(0, 1 - y\hat{y})$$

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

## A Perceptron Example

---

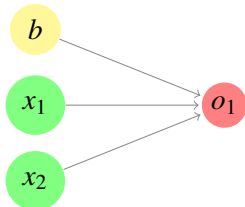
$$\mathbf{x} = (x_1, x_2), y = f(x_1, x_2)$$



## 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?

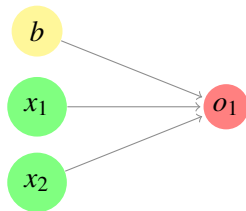
$x_1$	0	1	0	1
$x_2$	0	0	1	1
$y = x_1 \vee x_2$	0	1	1	1

## 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?

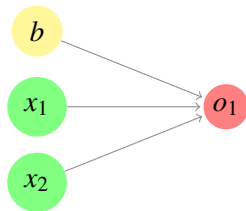
$x_1$	0	1	0	1
$x_2$	0	0	1	1
$y = x_1 \wedge x_2$	0	0	0	1

## 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

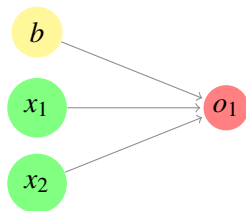
## A Perceptron Example

---

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.

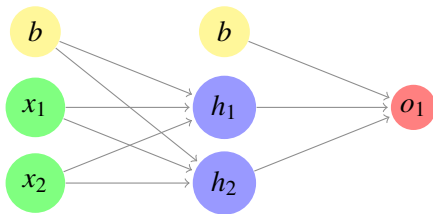
$x_1$		0	1	0	1
$x_2$		0	0	1	1
$x_1 \text{ XOR } x_2$		0	1	1	0

$$\text{XOR} = \text{AND} ( \text{OR} , \text{NAND} )$$

## 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



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

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

## General Expressiveness of Neural Networks

---

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

## Summary

---

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

## References

---

George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems (MCSS)*, 2(4):303–314, 1989.

Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.