

Neural Networks

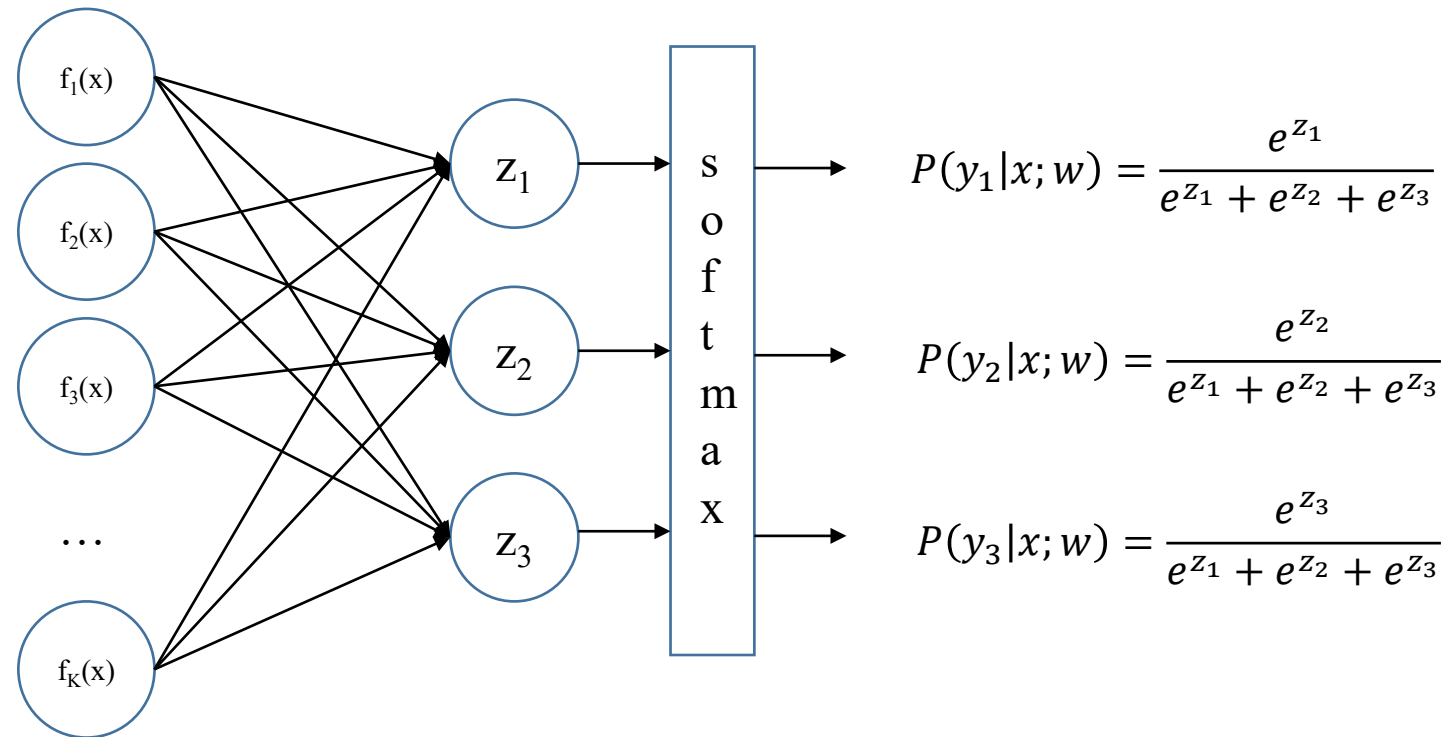
Kaiqi Zhao
The University of Auckland

Slides are partially based on the materials from Stanford University and UC Berkeley

Neural Networks

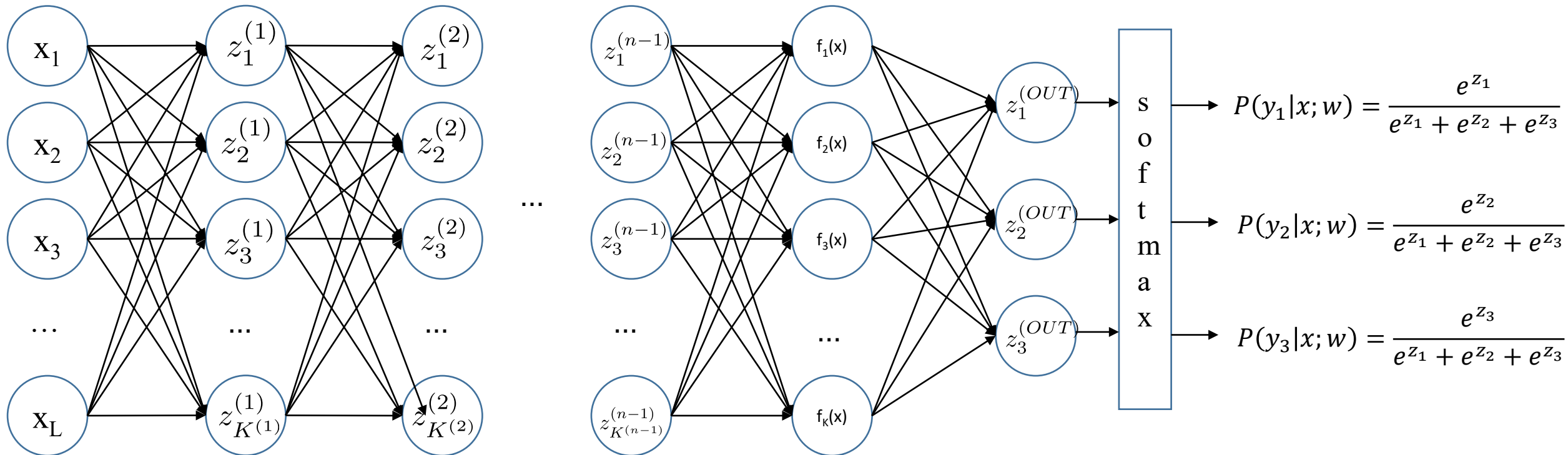
Multi-class Logistic Regression

- Multi-class Logistic Regression is **a special case of neural network**



Deep Neural Network = Also learn the features!

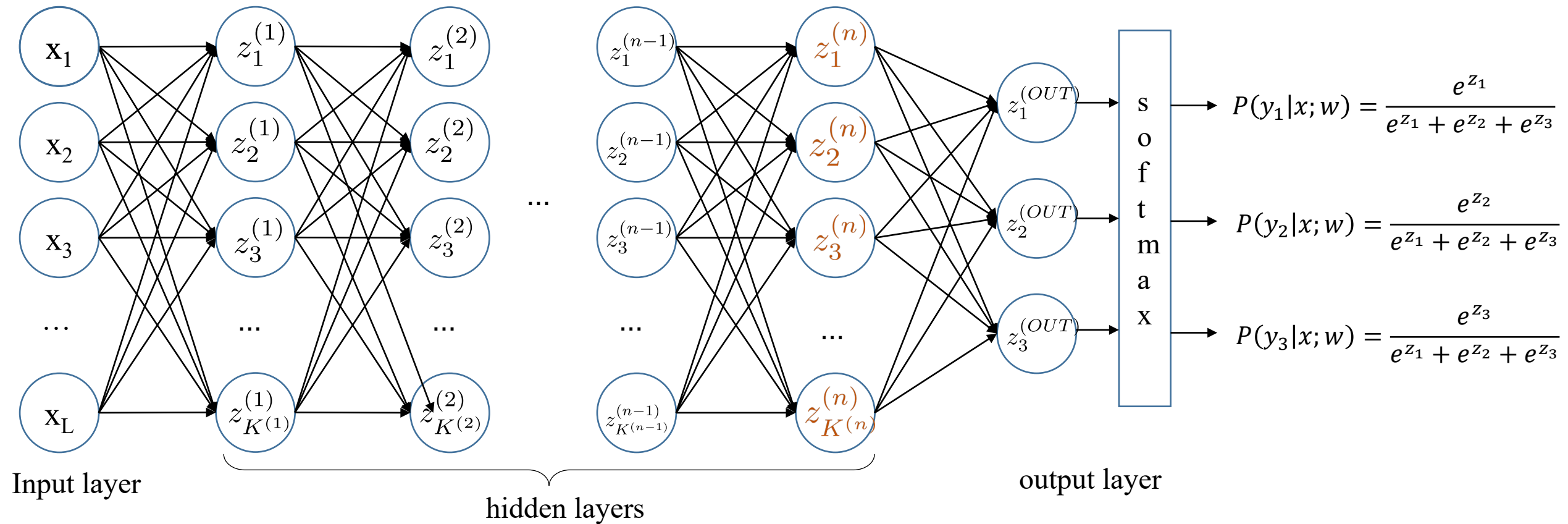
Neural Network = Also learn the features!



$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

g = nonlinear activation function

Neural Network = Also learn the features!



Intuitive Example

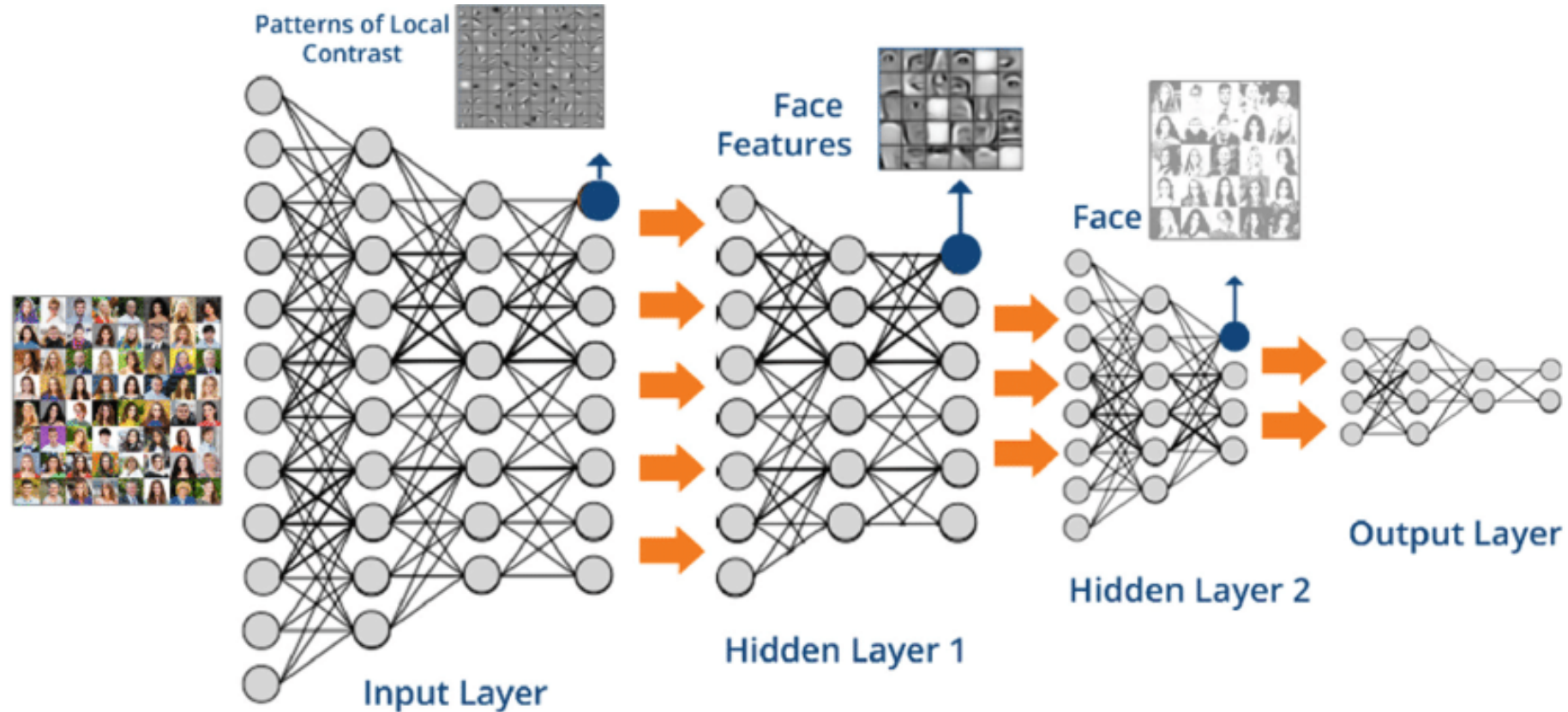
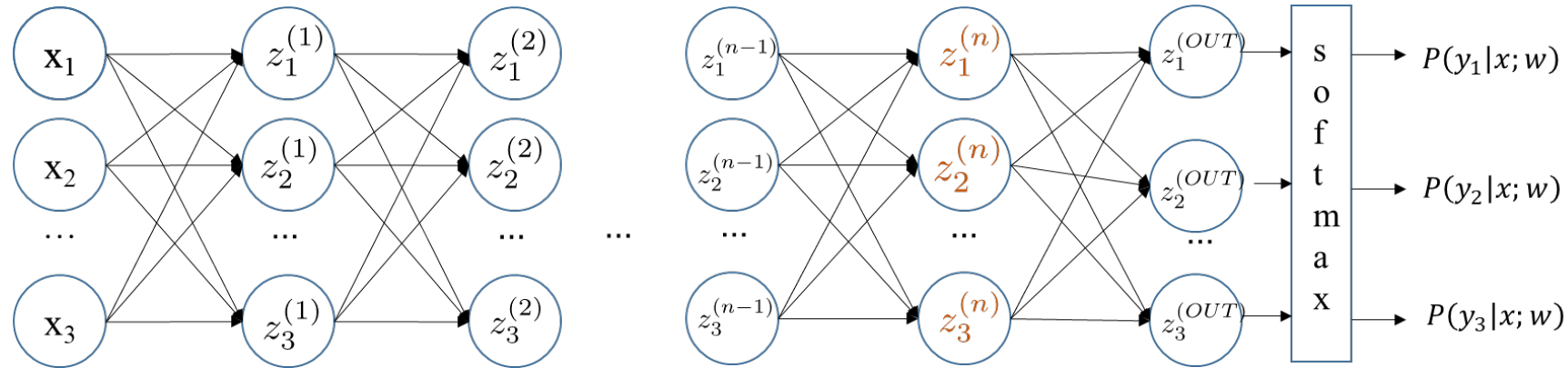


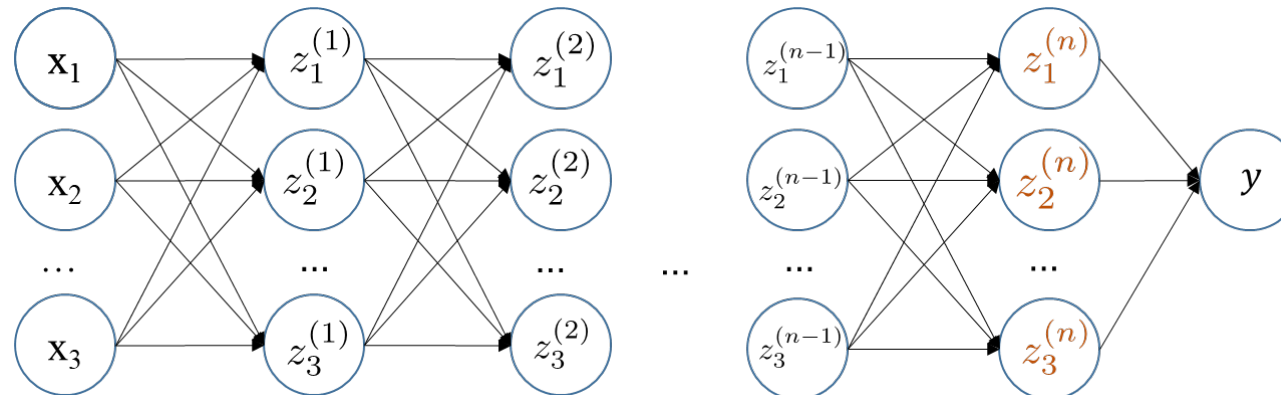
Image source: Grigsby, Scott. (2018). Artificial Intelligence for Advanced Human-Machine Symbiosis. 10.1007/978-3-319-91470-1_22.

Neural Network

Neural Networks for classification



Neural Networks for regression

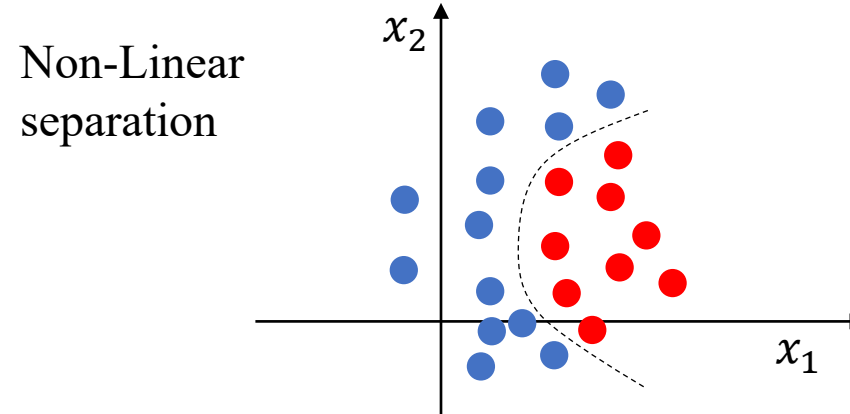
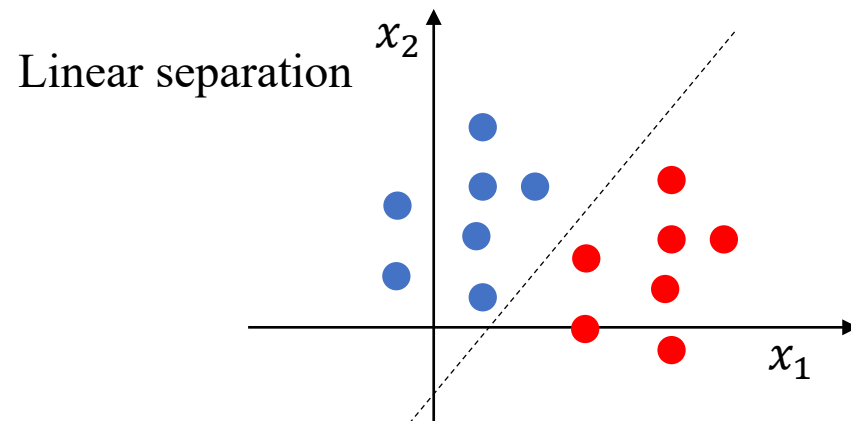
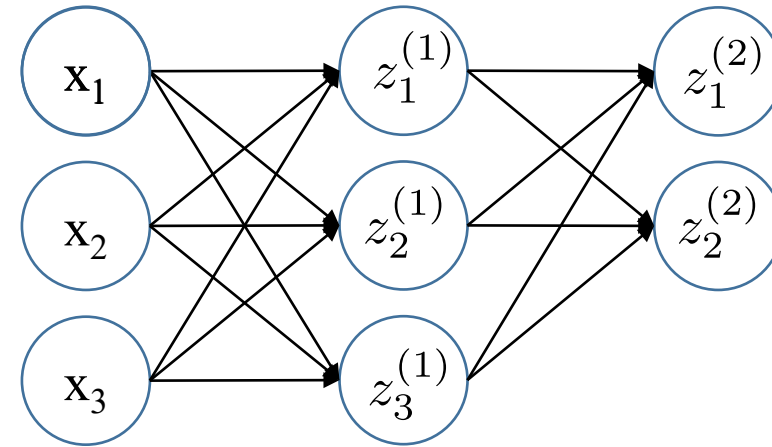


$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

g = nonlinear activation function

Why Non-linear Activation?

- How if we stack several linear layers?
 - For example two layers:
 - $z^{(1)} = w_1^T x$
 - $z^{(2)} = w_2^T z^{(1)}$
 - $z^{(2)} = w_2^T w_1^T x$
- It is equivalent to a single linear layer!!
- Some datasets are not linearly separable!!

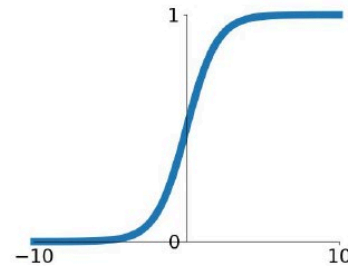


Activation Functions

- There are many choices for the activation function. We will cover them later

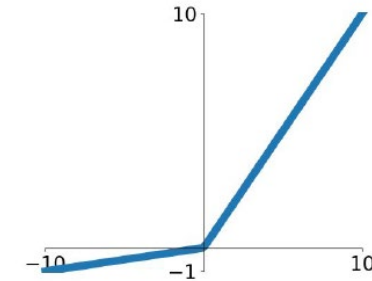
Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



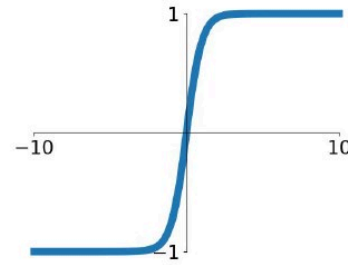
Leaky ReLU

$$\max(0.1x, x)$$



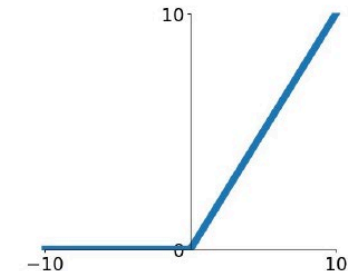
Tanh

$$\tanh(x)$$



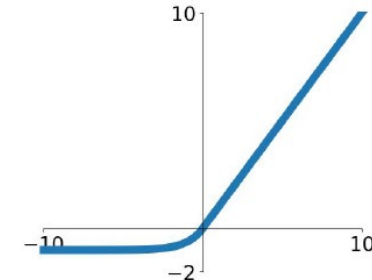
ReLU

$$\max(0, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Neural Networks Properties

- Theorem (Universal Function Approximators). A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.
- Practical considerations
 - Can be seen as learning the features
 - Large number of neurons
 - Danger for overfitting
 - (hence early stopping!)

Universal Function Approximation Theorem*

Hornik theorem 1: Whenever the activation function is *bounded and nonconstant*, then, for any finite measure μ , standard multilayer feedforward networks can approximate any function in $L^p(\mu)$ (the space of all functions on R^k such that $\int_{R^k} |f(x)|^p d\mu(x) < \infty$) arbitrarily well, provided that sufficiently many hidden units are available.

Hornik theorem 2: Whenever the activation function is *continuous, bounded and non-constant*, then, for arbitrary compact subsets $X \subseteq R^k$, standard multilayer feedforward networks can approximate any continuous function on X arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

- **In words:** Given any continuous function $f(x)$, if a 2-layer neural network has enough hidden units, then there is a choice of weights that allow it to closely approximate $f(x)$.

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"

Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"

Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

Universal Function Approximation Theorem*

Math. Control Signals Systems (1989) 2: 303–314

Mathematics of Control,
Signals, and Systems
© 1989 Springer-Verlag New York Inc.

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of n real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

Key words. Neural networks, Approximation, Completeness.

1. Introduction

A number of diverse application areas are concerned with the representation of general functions of an n -dimensional real variable, $x \in \mathbb{R}^n$, by finite linear combinations of the form

$$\sum_{j=1}^N \alpha_j \sigma(y_j^T x + \theta_j), \quad (1)$$

where $y_j \in \mathbb{R}^n$ and $\alpha_j, \theta_j \in \mathbb{R}$ are fixed. (y^T is the transpose of y so that $y^T x$ is the inner product of y and x .) Here the univariate function σ depends heavily on the context of the application. Our major concern is with so-called sigmoidal σ 's:

$$\sigma(t) \rightarrow \begin{cases} 1 & \text{as } t \rightarrow +\infty, \\ 0 & \text{as } t \rightarrow -\infty. \end{cases}$$

Such functions arise naturally in neural network theory as the activation function of a neural node (or unit as is becoming the preferred term) [L1], [RHM]. The main result of this paper is a demonstration of the fact that sums of the form (1) are dense in the space of continuous functions on the unit cube if σ is any continuous sigmoidal

* Date received: October 21, 1988. Date revised: February 17, 1989. This research was supported in part by NSF Grant DCR-8619103, ONR Contract N000-86-G-0202 and DOE Grant DE-FG02-85ER25001.

† Center for Supercomputing Research and Development and Department of Electrical and Computer Engineering, University of Illinois, Urbana, Illinois 61801, U.S.A.

Neural Networks, Vol. 4, pp. 251–257, 1991
Printed in the USA. All rights reserved.

(0893-6080/91 \$3.00 + .00
Copyright © 1991 Pergamon Press plc

ORIGINAL CONTRIBUTION

Approximation Capabilities of Multilayer Feedforward Networks

KURT HORNİK

Technische Universität Wien, Vienna, Austria

(Received 30 January 1990; revised and accepted 25 October 1990)

Abstract—We show that standard multilayer feedforward networks with as few as a single hidden layer and arbitrary bounded and nonconstant activation function are universal approximators with respect to $L^p(\mu)$ performance criteria, for arbitrary finite input environment measures μ , provided only that sufficiently many hidden units are available. If the activation function is continuous, bounded and nonconstant, then continuous mappings can be learned uniformly over compact input sets. We also give very general conditions ensuring that networks with sufficiently smooth activation functions are capable of arbitrarily accurate approximation to a function and its derivatives.

Keywords—Multilayer feedforward networks, Activation function, Universal approximation capabilities, Input environment measure, $L^p(\mu)$ approximation, Uniform approximation, Sobolev spaces, Smooth approximation.

1. INTRODUCTION

The approximation capabilities of neural network architectures have recently been investigated by many authors, including Carroll and Dickinson (1989), Cybenko (1989), Funahashi (1989), Gallant and White (1988), Hecht-Nielsen (1989), Hornik, Stinchcombe, and White (1989, 1990), Irie and Miyake (1988), Lapedes and Farber (1988), Stinchcombe and White (1989, 1990). (This list is by no means complete.)

If we think of the network architecture as a rule for computing values at l output units given values at k input units, hence implementing a class of mappings from \mathbb{R}^k to \mathbb{R}^l , we can ask how well arbitrary mappings from \mathbb{R}^k to \mathbb{R}^l can be approximated by the network, in particular, if as many hidden units as required for internal representation and computation may be employed.

How to measure the accuracy of approximation depends on how we measure closeness between functions, which in turn varies significantly with the specific problem to be dealt with. In many applications, it is necessary to have the network perform simultaneously well on all input samples taken from some compact input set X in \mathbb{R}^k . In this case, closeness is

measured by the uniform distance between functions on X , that is,

$$\rho_{\infty}(f, g) = \sup_{x \in X} |f(x) - g(x)|.$$

In other applications, we think of the inputs as random variables and are interested in the average performance where the average is taken with respect to the input environment measure μ , where $\mu(\mathbb{R}^k) < \infty$. In this case, closeness is measured by the $L^p(\mu)$ distances

$$\rho_p(f, g) = \left[\int_X |f(x) - g(x)|^p d\mu(x) \right]^{1/p},$$

$1 \leq p < \infty$, the most popular choice being $p = 2$, corresponding to mean square error.

Of course, there are many more ways of measuring closeness of functions. In particular, in many applications, it is also necessary that the derivatives of the approximating function implemented by the network closely resemble those of the function to be approximated, up to some order. This issue was first taken up in Hornik et al. (1990), who discuss the sources of need of smooth functional approximation in more detail. Typical examples arise in robotics (learning of smooth movements) and signal processing (analysis of chaotic time series); for a recent application to problems of nonparametric inference in statistics and econometrics, see Gallant and White (1989).

All papers establishing certain approximation ca-

Requests for reprints should be sent to Kurt Hornik, Institut für Statistik und Wahrscheinlichkeitstheorie, Technische Universität Wien, Wiedner Hauptstraße 8-10/107, A-1040 Wien, Austria.

MULTILAYER FEEDFORWARD NETWORKS WITH NON-POLYNOMIAL ACTIVATION FUNCTIONS CAN APPROXIMATE ANY FUNCTION

by

Moshe Leshno
Faculty of Management
Tel Aviv University
Tel Aviv, Israel 69978

and

Shimon Schocken
Leonard N. Stern School of Business
New York University
New York, NY 10003

September 1991

Center for Research on Information Systems
Information Systems Department
Leonard N. Stern School of Business
New York University

Working Paper Series

STERN IS-91-26

Appeared previously as Working Paper No. 21/91 at The Israel Institute Of Business Research

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

Training a Neural Network

- Training a neural network for classification is just like logistic regression:
 - Maximize the log-likelihood, or
 - Minimize the cross-entropy loss
- Training a neural network for regression is to minimize the sum of square error
- Run gradient ascent/descent! → we need to compute the derivatives

Computing the derivatives

- Derivatives tables:

$$\frac{d}{dx}(a) = 0$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(au) = a \frac{du}{dx}$$

$$\frac{d}{dx}(u + v - w) = \frac{du}{dx} + \frac{dv}{dx} - \frac{dw}{dx}$$

$$\frac{d}{dx}(uv) = u \frac{dv}{dx} + v \frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{u}{v}\right) = \frac{1}{v} \frac{du}{dx} - \frac{u}{v^2} \frac{dv}{dx}$$

$$\frac{d}{dx}(u^n) = nu^{n-1} \frac{du}{dx}$$

$$\frac{d}{dx}(\sqrt{u}) = \frac{1}{2\sqrt{u}} \frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u}\right) = -\frac{1}{u^2} \frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u^n}\right) = -\frac{n}{u^{n+1}} \frac{du}{dx}$$

$$\frac{d}{dx}[f(u)] = \frac{d}{du}[f(u)] \frac{du}{dx}$$

$$\frac{d}{dx}[\ln u] = \frac{d}{dx}[\log_e u] = \frac{1}{u} \frac{du}{dx}$$

$$\frac{d}{dx}[\log_a u] = \log_a e \frac{1}{u} \frac{du}{dx}$$

$$\frac{d}{dx}e^u = e^u \frac{du}{dx}$$

$$\frac{d}{dx}a^u = a^u \ln a \frac{du}{dx}$$

$$\frac{d}{dx}(u^v) = vu^{v-1} \frac{du}{dx} + \ln u \ u^v \frac{dv}{dx}$$

$$\frac{d}{dx} \sin u = \cos u \frac{du}{dx}$$

$$\frac{d}{dx} \cos u = -\sin u \frac{du}{dx}$$

$$\frac{d}{dx} \tan u = \sec^2 u \frac{du}{dx}$$

$$\frac{d}{dx} \cot u = -\csc^2 u \frac{du}{dx}$$

$$\frac{d}{dx} \sec u = \sec u \tan u \frac{du}{dx}$$

$$\frac{d}{dx} \csc u = -\csc u \cot u \frac{du}{dx}$$

Computing the derivatives

- But neural net f is never one of those?
- No problem: CHAIN RULE:

If $f(x) = g(h(x))$

Then
$$\frac{df(x)}{dx} = \frac{dg}{dh} \cdot \frac{dh(x)}{dx}$$

→ Derivatives can be computed by following well-defined procedures

Example

- The cross-entropy for binary classification problem:

$$CE(w) = -y \ln \hat{y} - (1 - y) \ln(1 - \hat{y}), \quad \hat{y} = p(y = +1|x, w) = \frac{1}{1 + e^{-wx}}$$

- By chain rule:

$$\frac{dCE(w)}{dw} = \frac{dCE(w)}{d\hat{y}} \cdot \frac{d\hat{y}}{dw}$$

$$\frac{dCE(w)}{d\hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} = \frac{\hat{y}-y}{\hat{y}(1-\hat{y})} \quad \frac{d\hat{y}}{dw} = ?$$

Let $h = -wx$, then $\hat{y} = \frac{1}{1+e^h}$

$$\frac{d\hat{y}}{dw} = \frac{d\hat{y}}{dh} \cdot \frac{dh}{dw} = -\frac{e^h}{(1+e^h)^2} \cdot -x = \hat{y}(1-\hat{y})x$$

$$\frac{dCE(w)}{dw} = \frac{dCE(w)}{d\hat{y}} \cdot \frac{d\hat{y}}{dw}$$

$$= \frac{\hat{y}-y}{\hat{y}(1-\hat{y})} \cdot \hat{y}(1-\hat{y})x$$

$$= (\hat{y}-y)x$$

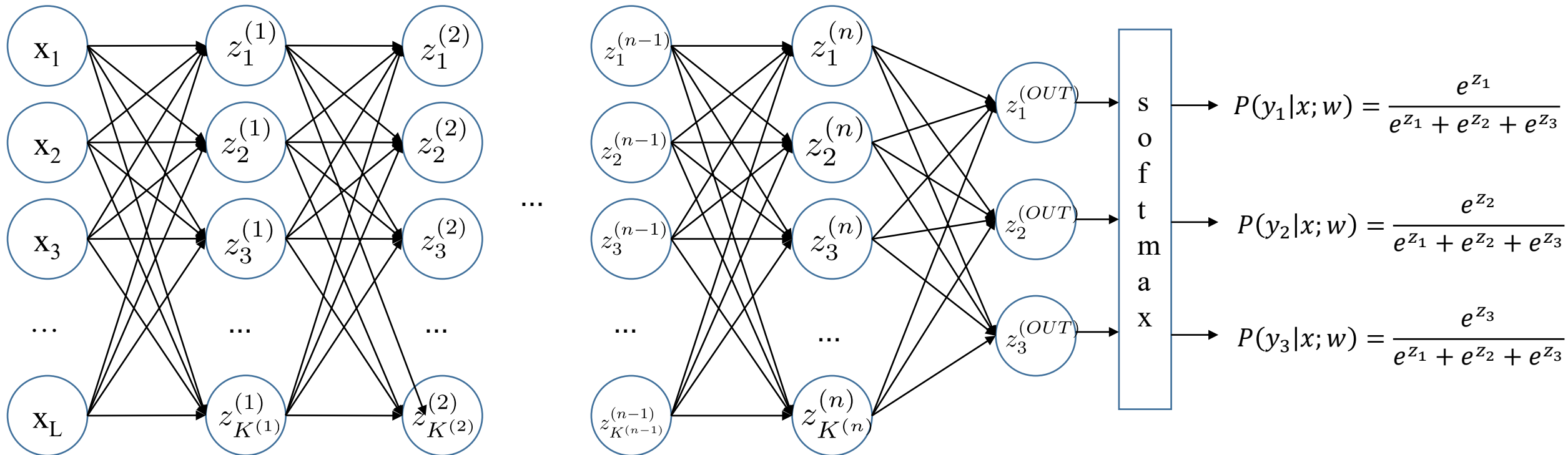
Gradient descent:

$$w \leftarrow w - \alpha \cdot (\hat{y} - y)x$$

Automatic Differentiation

- Automatic differentiation software
 - e.g. PyTorch, TensorFlow
 - Only need to program the function $g(x, y, w)$
 - Can automatically compute all derivatives w.r.t. all entries in w
 - This is typically done by caching info during forward computation pass of your NN f , and then doing a backward pass = “backpropagation”
 - Autodiff / Backpropagation can often be done at computational cost comparable to the forward pass
- You only need to know this exists

Training a Neural Network



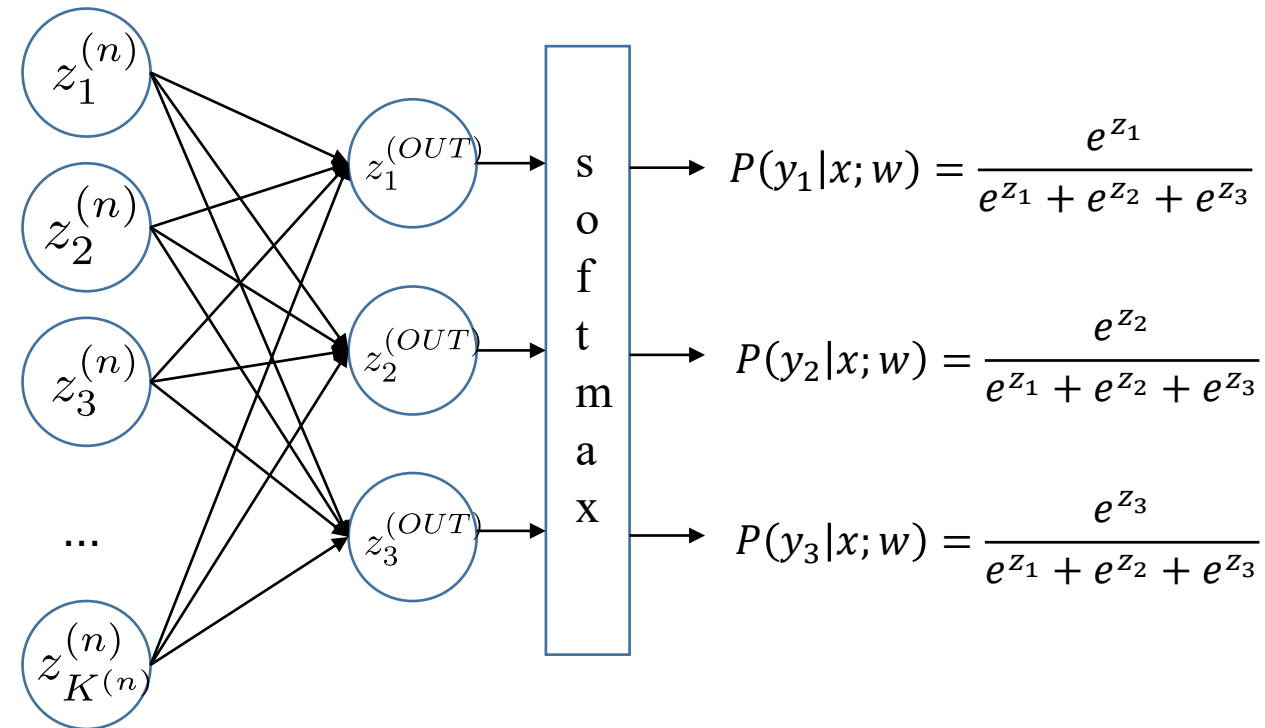
$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

g = nonlinear activation function

Training a Neural Network

Key words:

- Forward
- Backwards
- Gradient
- Backprop



g = nonlinear activation function

Training a Neural Network

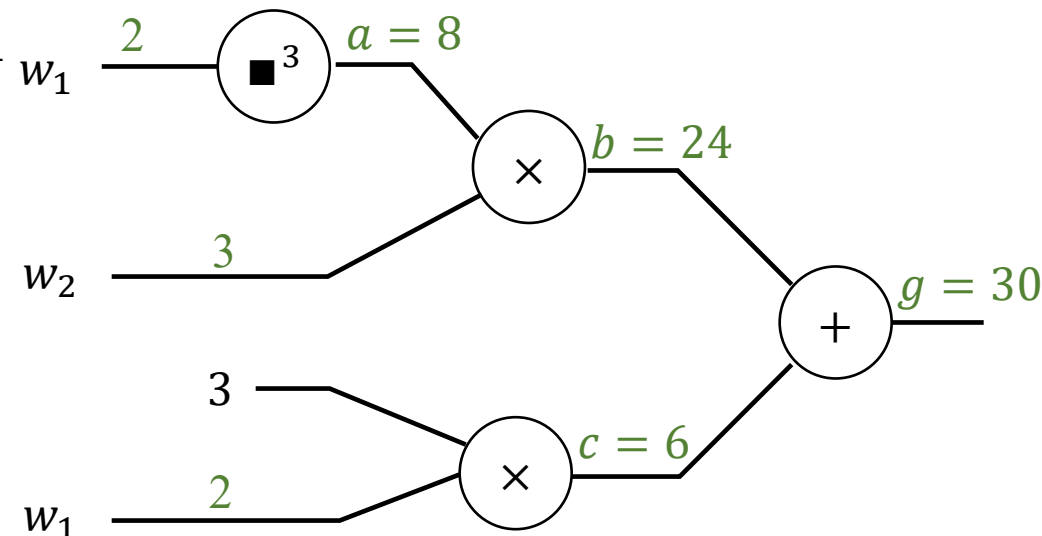
- Let's use a simple example to learn how to train a Neural Network.
- Basically, the loss function (e.g., cross-entropy) of a Neural Network can be considered as a complex function that maps your parameters w to the loss.

- **A simple (loss) function:** $g(w) = w_1^3 w_2 + 3w_1$

- A computational graph could be constructed from the objective function

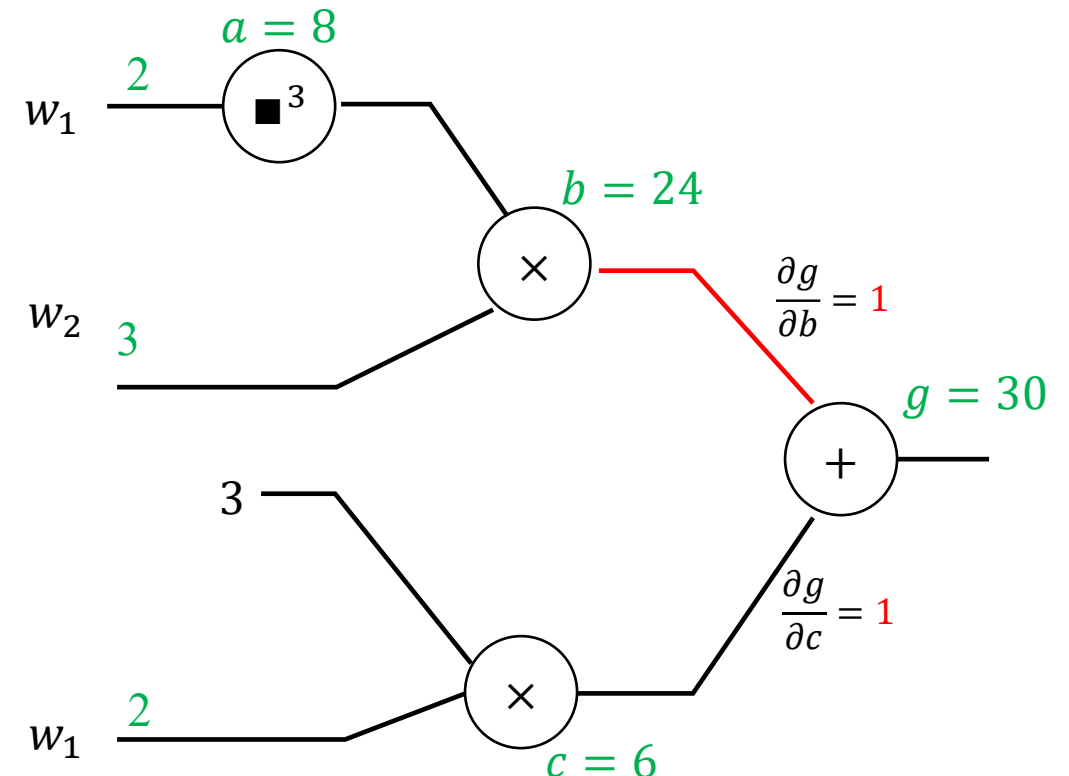
- **Forward step** – substitute your input and get values for each node in the computational graph

- Calculate the gradient backward by **backpropagation**



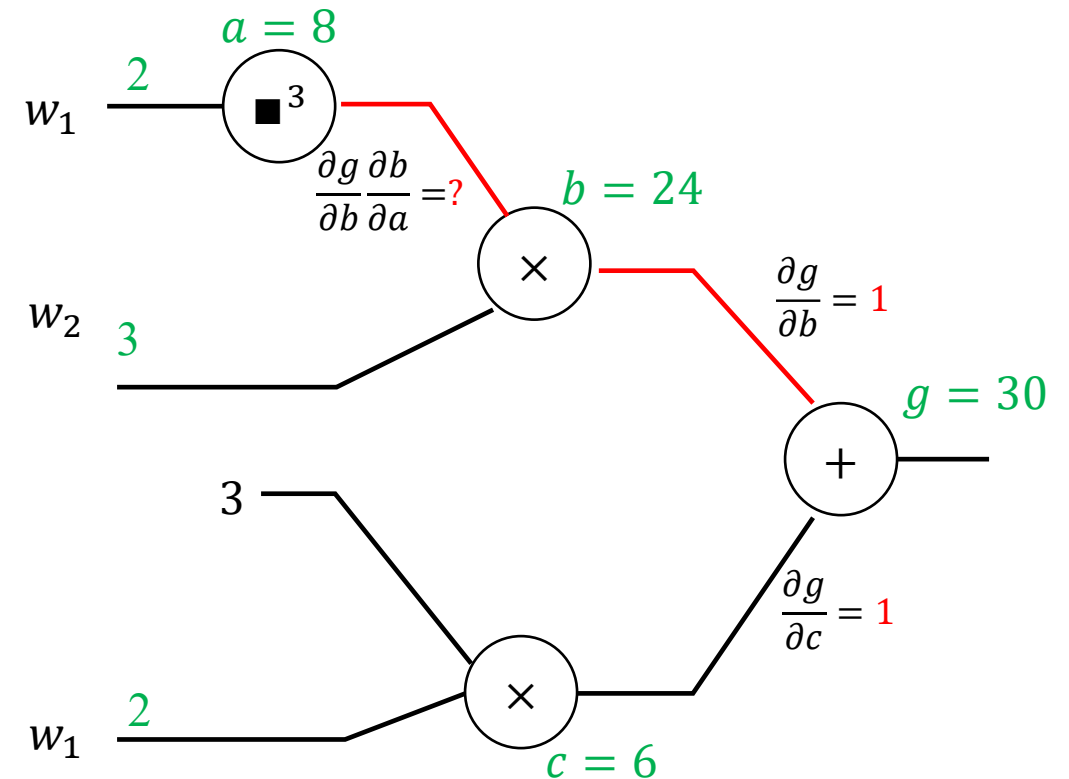
Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.
- $g = b + c$
 - $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.
- $g = b + c$
 - $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$
- $b = a \times w_2$
 - $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = ? ? ? ? ?$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

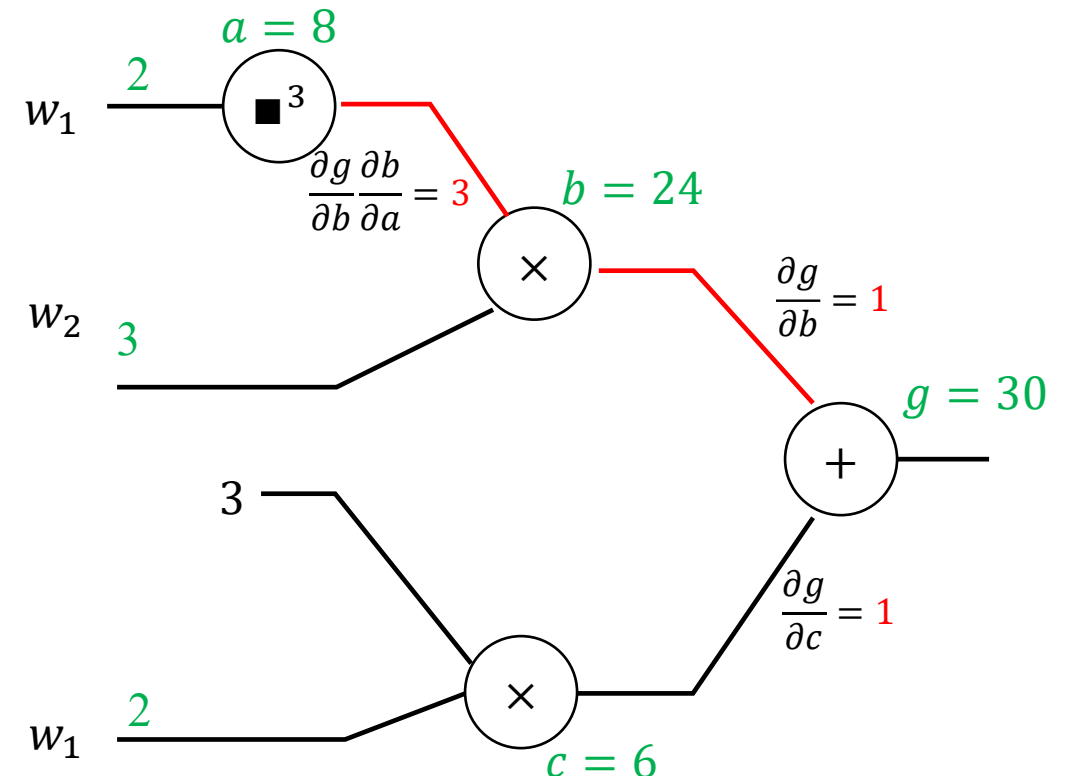
- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.

- $g = b + c$

- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.

- $g = b + c$

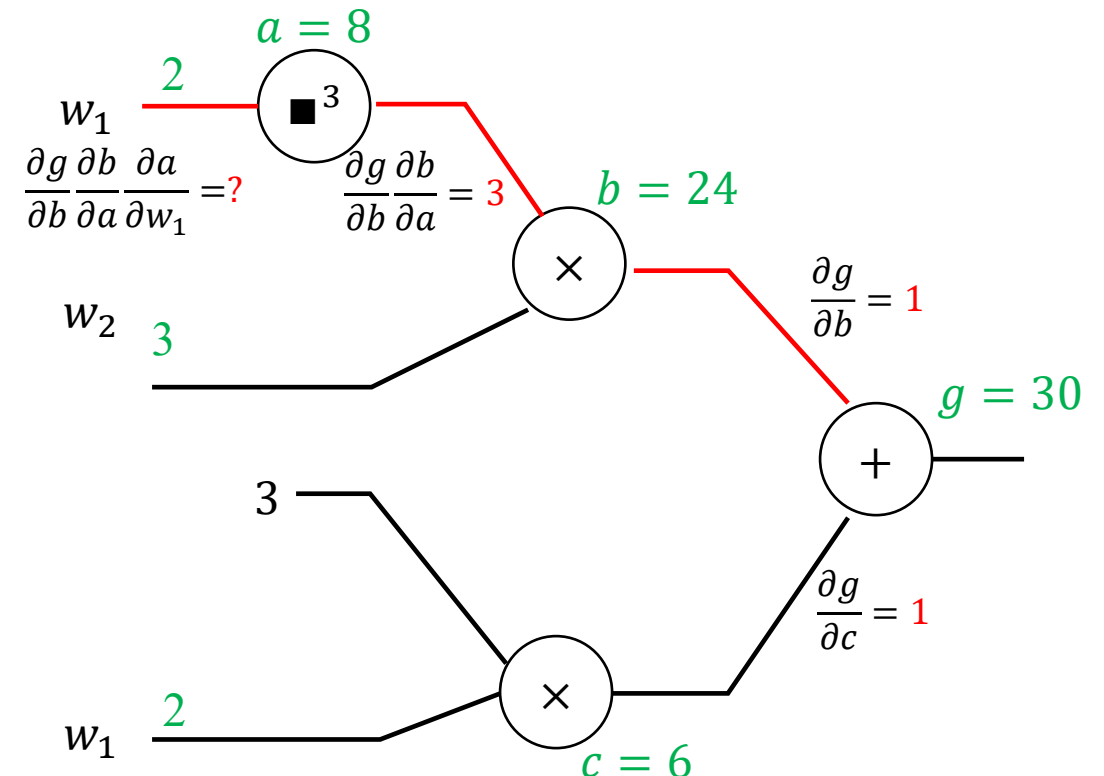
- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$

- $a = w_1^3$

- $\frac{\partial g}{\partial w_1} = \text{?????}$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.

- $g = b + c$

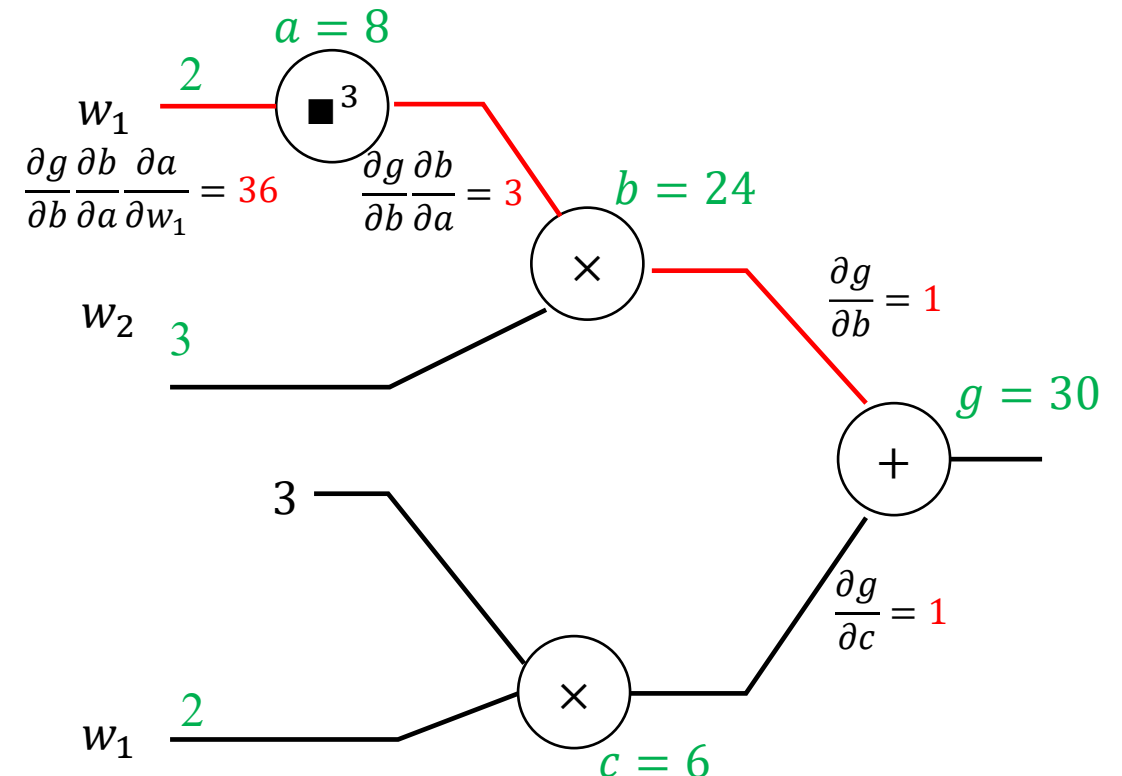
- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$

- $a = w_1^3$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial a} \frac{\partial a}{\partial w_1} = 3 \cdot 3w_1^2 = 36$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.

- $g = b + c$

- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

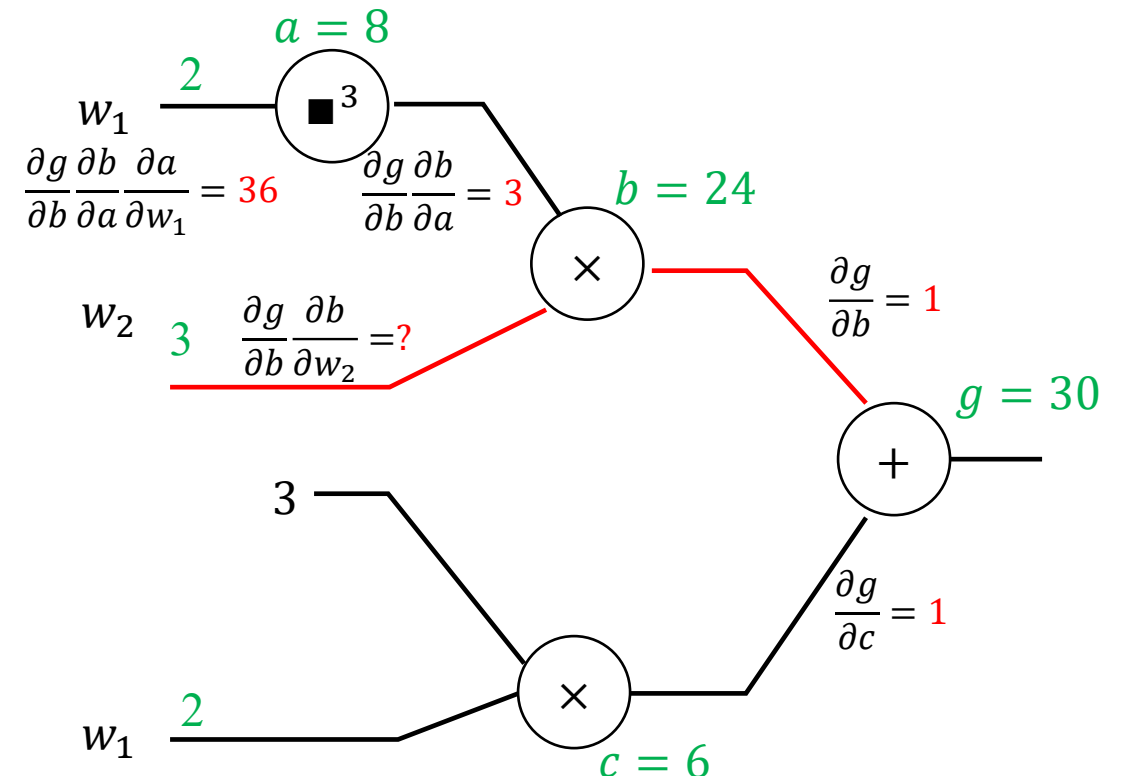
- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$

- $a = w_1^3$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial a} \frac{\partial a}{\partial w_1} = 3 \cdot 3w_1^2 = 36$

- $\frac{\partial g}{\partial w_2} = ???$ Hint: $b = a \times w_2$ may be useful.



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.
 - Can use derivative chain rule to compute $\partial g / \partial w_1$ and $\partial g / \partial w_2$.

- $g = b + c$

- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

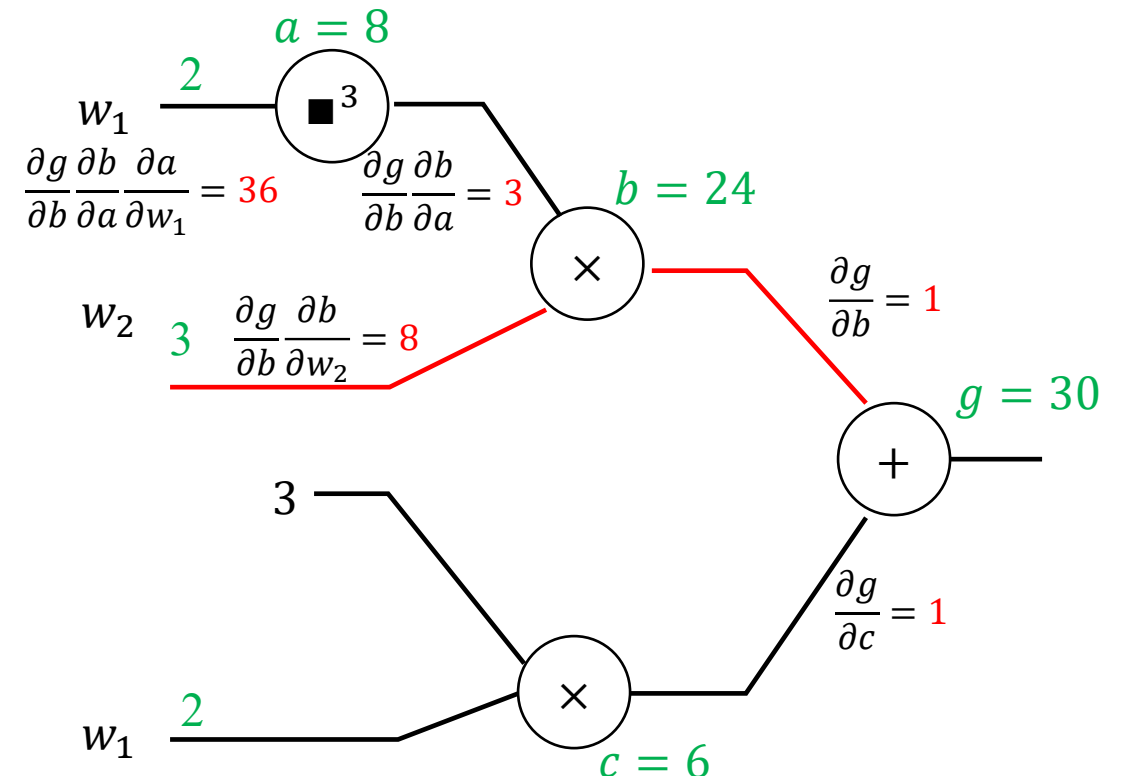
- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$

- $\frac{\partial g}{\partial w_2} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial w_2} = 1 \frac{\partial b}{\partial w_2} = 1 \cdot 8 = 8$

- $a = w_1^3$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial a} \frac{\partial a}{\partial w_1} = 3 \cdot 3w_1^2 = 36$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.

- $g = b + c$

- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$

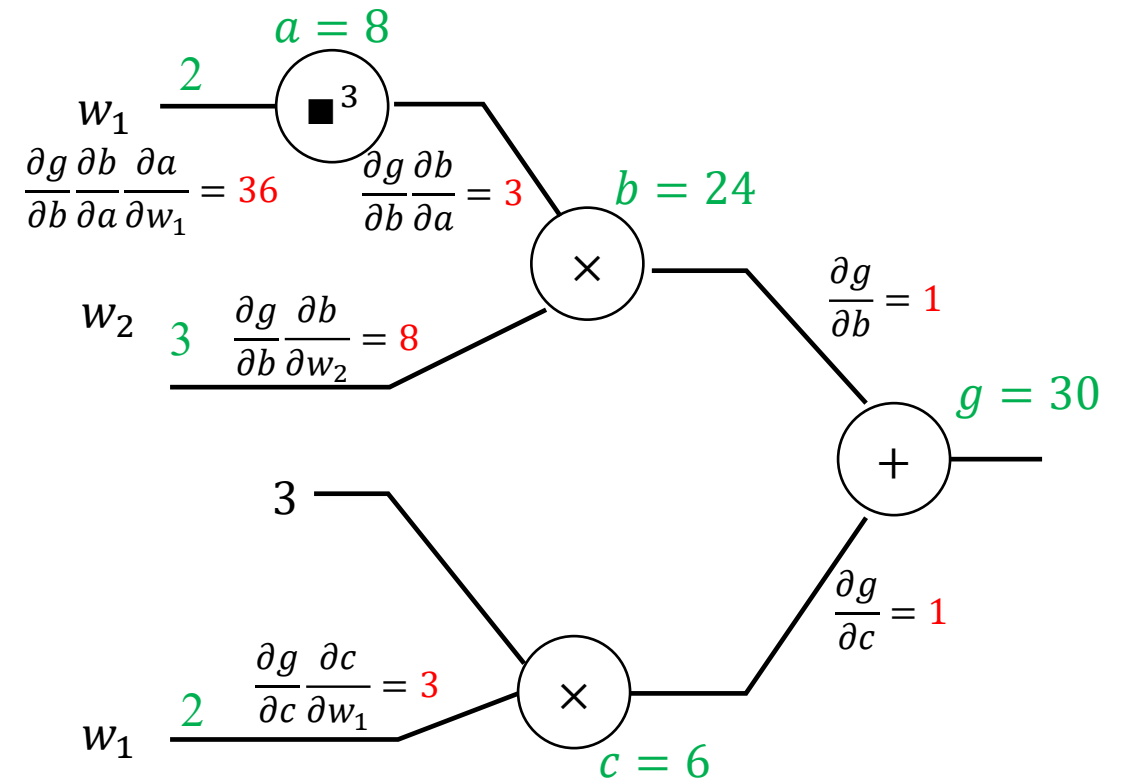
- $\frac{\partial g}{\partial w_2} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial w_2} = 1 \frac{\partial b}{\partial w_2} = 1 \cdot 8 = 8$

- $a = w_1^3$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial a} \frac{\partial a}{\partial w_1} = 3 \cdot 3w_1^2 = 36$

- $c = 3w_1$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial c} \frac{\partial c}{\partial w_1} = 1 \cdot 3 = 3$



Back Propagation: $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$

- Suppose we have $g(\mathbf{w}) = w_1^3 w_2 + 3w_1$ and want the gradient at $\mathbf{w} = [2, 3]$
- Think of the function as a composition of many functions.

- $g = b + c$

- $\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$

- $b = a \times w_2$

- $\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial a} = 1 \frac{\partial b}{\partial a} = 1 \cdot 3 = 3$

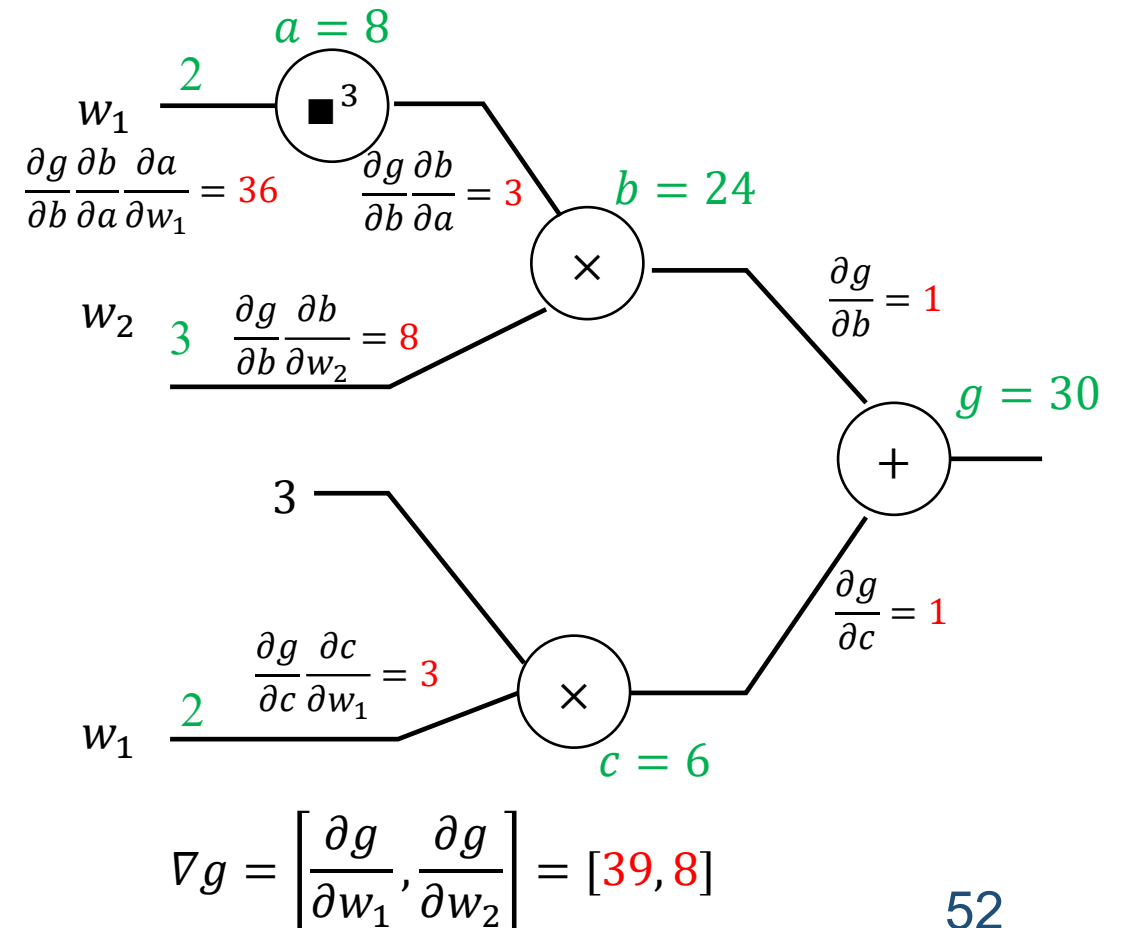
- $\frac{\partial g}{\partial w_2} = \frac{\partial g}{\partial b} \frac{\partial b}{\partial w_2} = 1 \frac{\partial b}{\partial w_2} = 1 \cdot 8 = 8$

- $a = w_1^3$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial a} \frac{\partial a}{\partial w_1} = 3 \cdot 3w_1^2 = 36$

- $c = 3w_1$

- $\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial c} \frac{\partial c}{\partial w_1} = 1 \cdot 3 = 3$

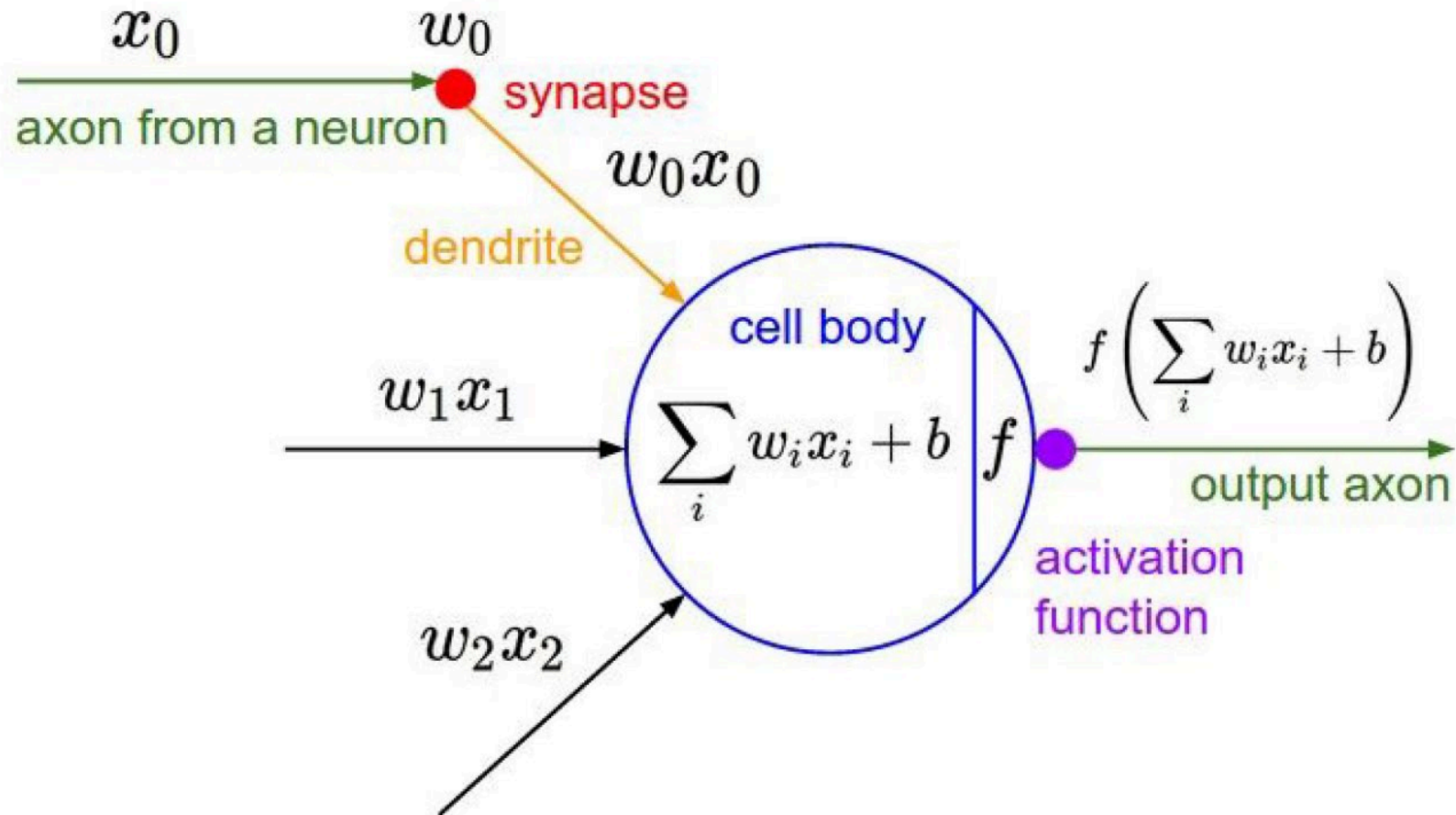


Summary of Key Ideas

- Optimize probability of label given input $\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$
- Continuous optimization
 - Gradient ascent:
 - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
 - Take step in the gradient direction
 - Repeat (until held-out data accuracy starts to drop = “early stopping”)
- Deep neural nets
 - Last layer = still logistic regression
 - Now also many more layers before this last layer
 - = computing the features
 - → the features are learned rather than hand-designed
 - Universal function approximation theorem
 - If the neural net is large enough, then it can represent any continuous mapping from input to output with arbitrary accuracy
 - But remember: need to avoid overfitting / memorizing the training data → early stopping!

Activation Functions

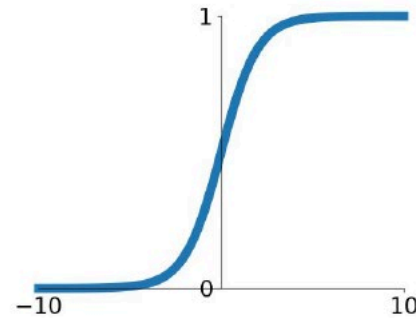
Activation Functions



Activation Functions

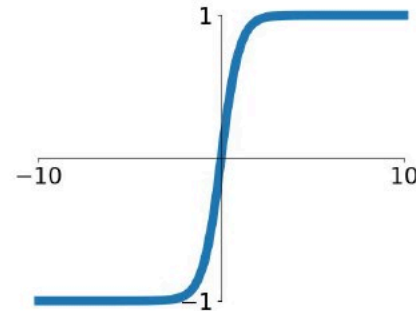
Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



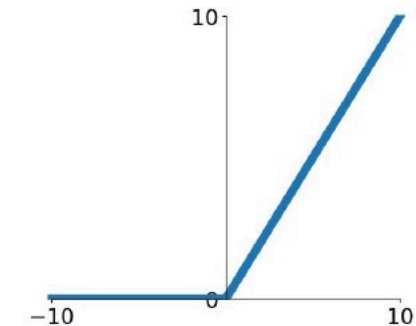
Tanh

$$\tanh(x)$$



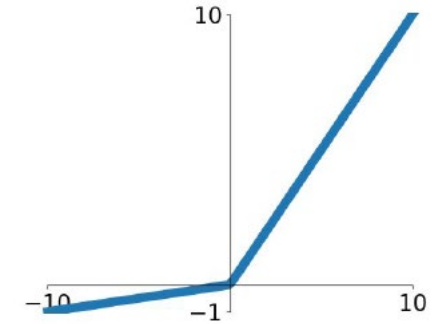
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$



ELU

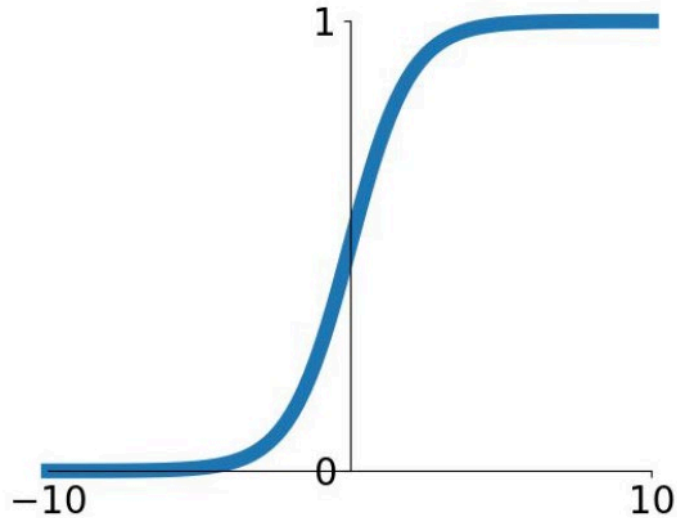
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

Reference:

[Activation Functions — ML Glossary documentation \(ml-cheatsheet.readthedocs.io\)](https://ml-cheatsheet.readthedocs.io/)

Activation Functions: Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



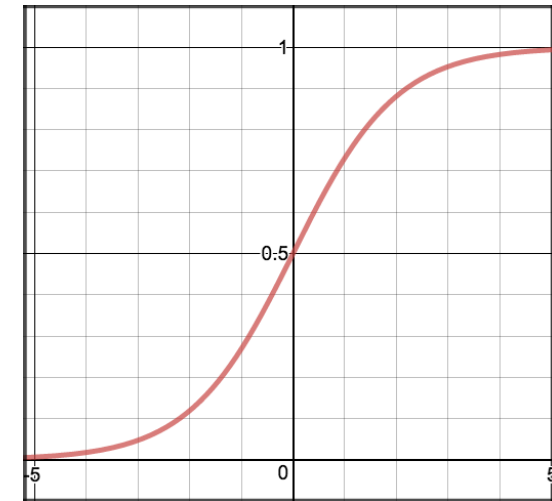
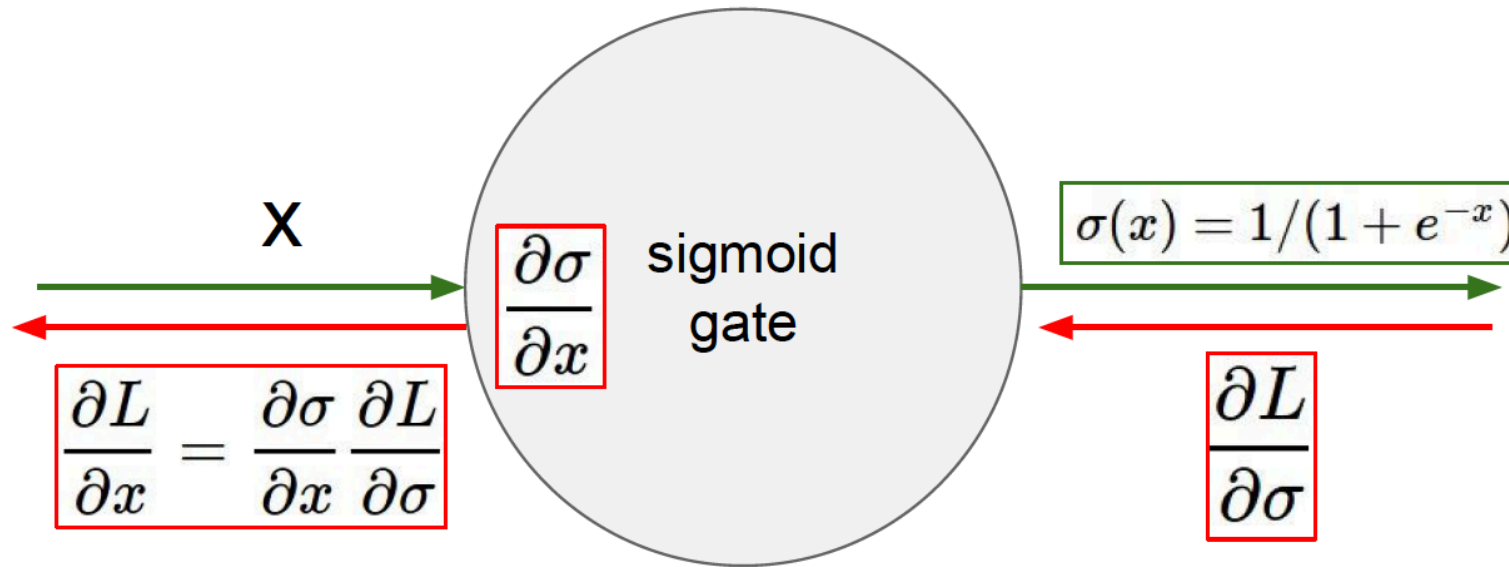
Sigmoid

- Squashes numbers to range $[0,1]$
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

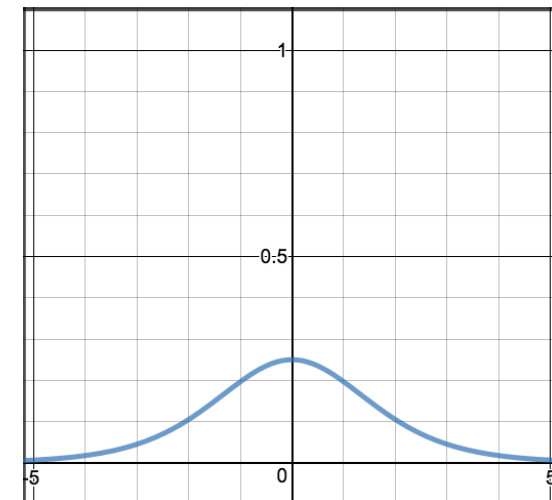
Three problems:

1. Saturated neurons “kill” the gradients

Activation Functions: Sigmoid



$\sigma(x)$

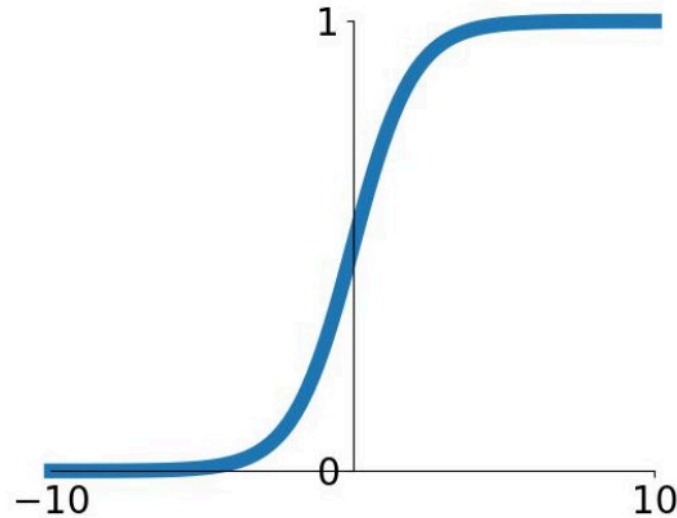


$\frac{\partial \sigma}{\partial x}(x)$

- What happens when $x < -5$?
- What happens when $x = 0$?
- What happens when $x > 5$?

Activation Functions: Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Sigmoid

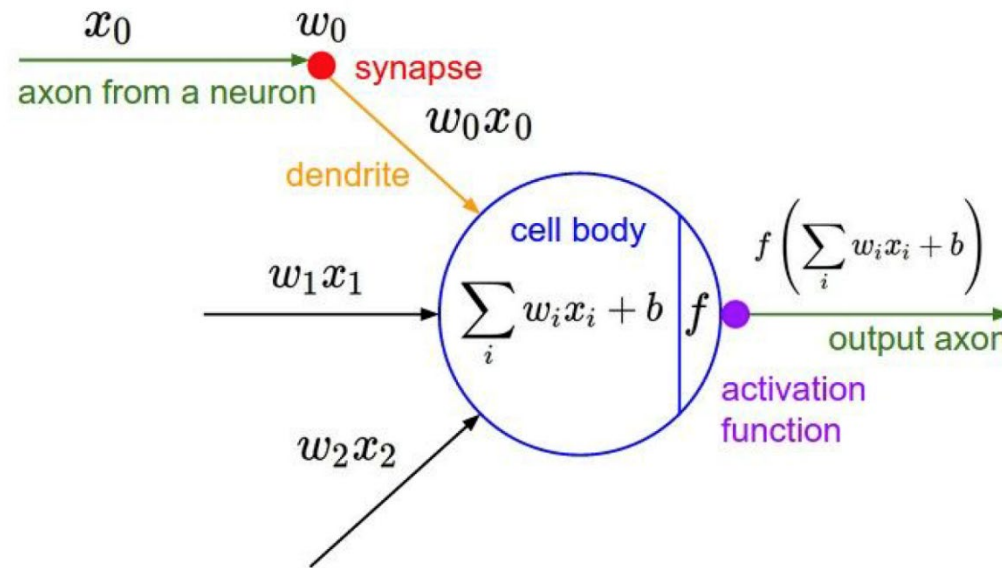
- Squashes numbers to range $[0,1]$
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

Three problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered

Activation Functions: Sigmoid

- Consider what happens when the input to a neuron (x) is always positive:



$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on \mathbf{w} ?

Activation Functions: Sigmoid

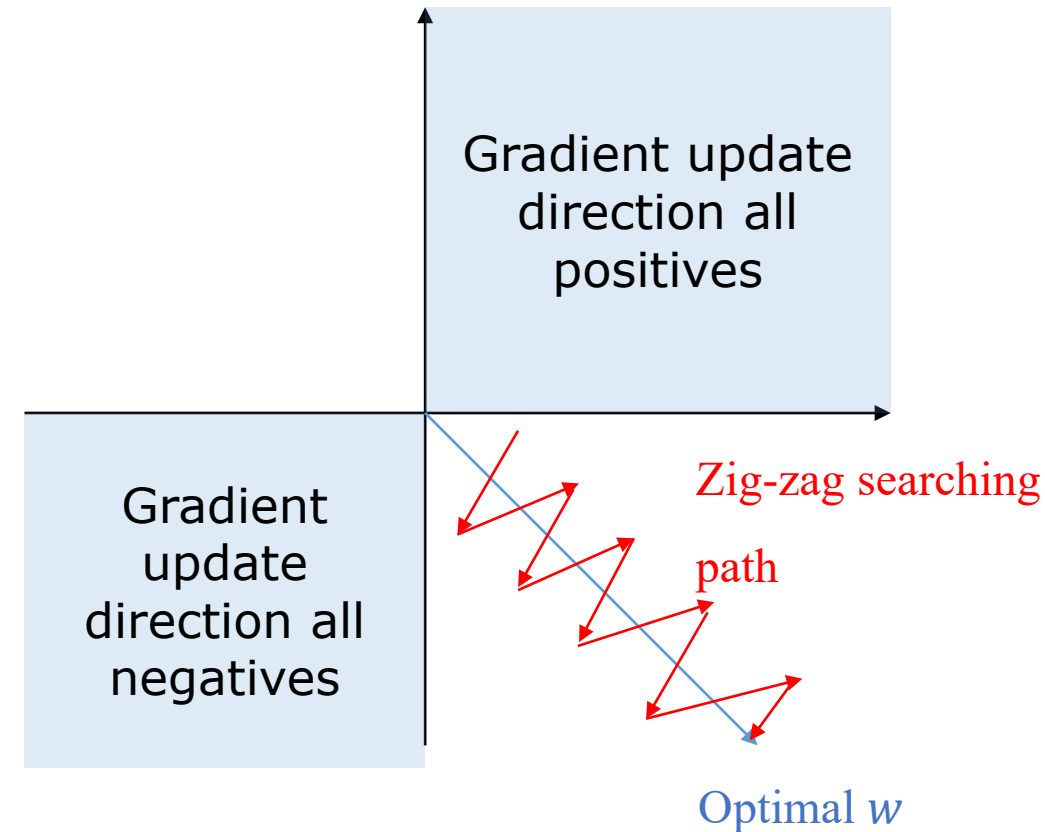
- Consider what happens when the input to a neuron (x) is always positive:

$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on \mathbf{w} ?

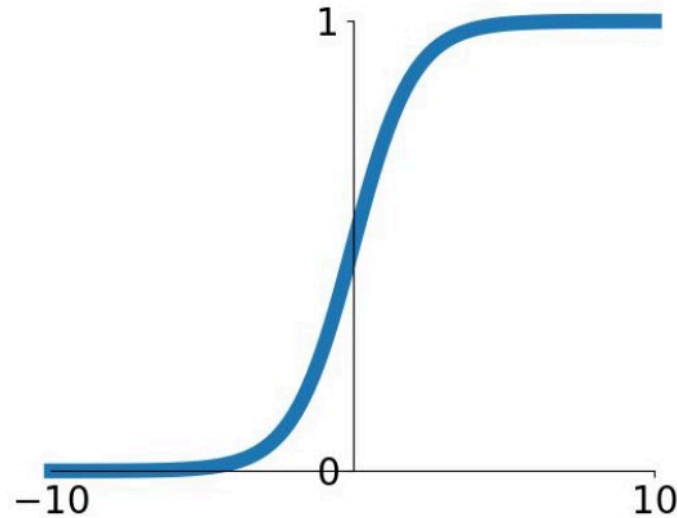
Always all positive or all negative

A zero-centered activation is preferred to avoid the zig-zag path



Activation Functions: Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



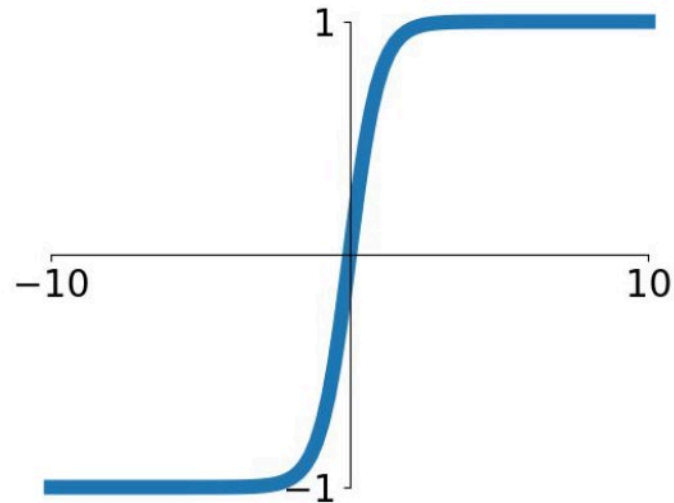
Sigmoid

- Squashes numbers to range $[0,1]$
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

Three problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit computationally expensive

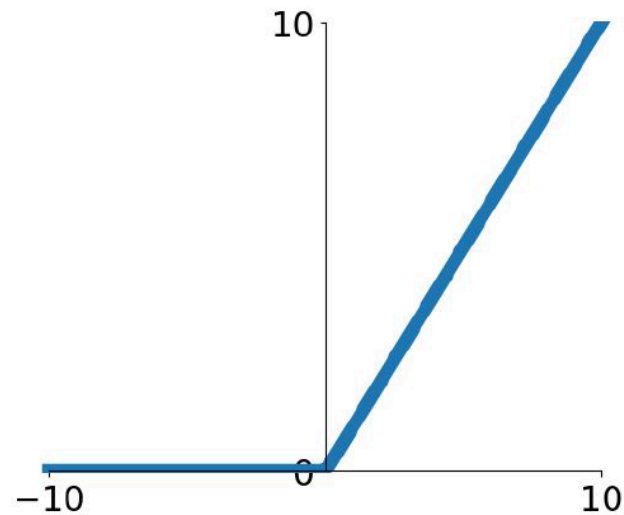
Activation Functions: Tanh



Tanh

- Squashes numbers to range $[-1,1]$
- zero centered (nice)
- still kills gradients when saturated :(

Activation Functions: ReLU

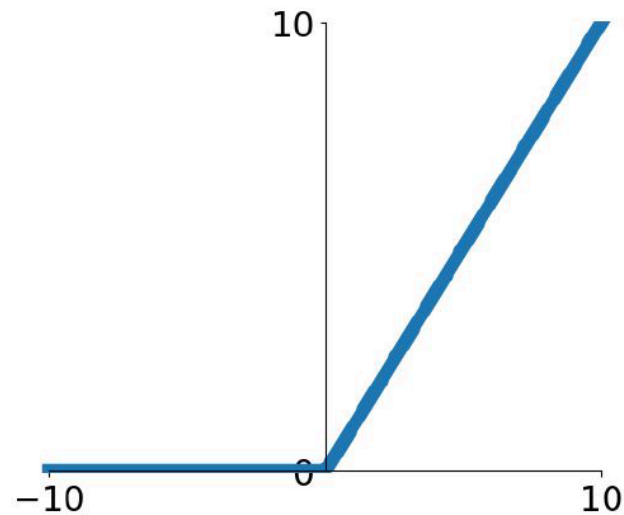


ReLU
(Rectified Linear Unit)

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

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

Activation Functions: ReLU



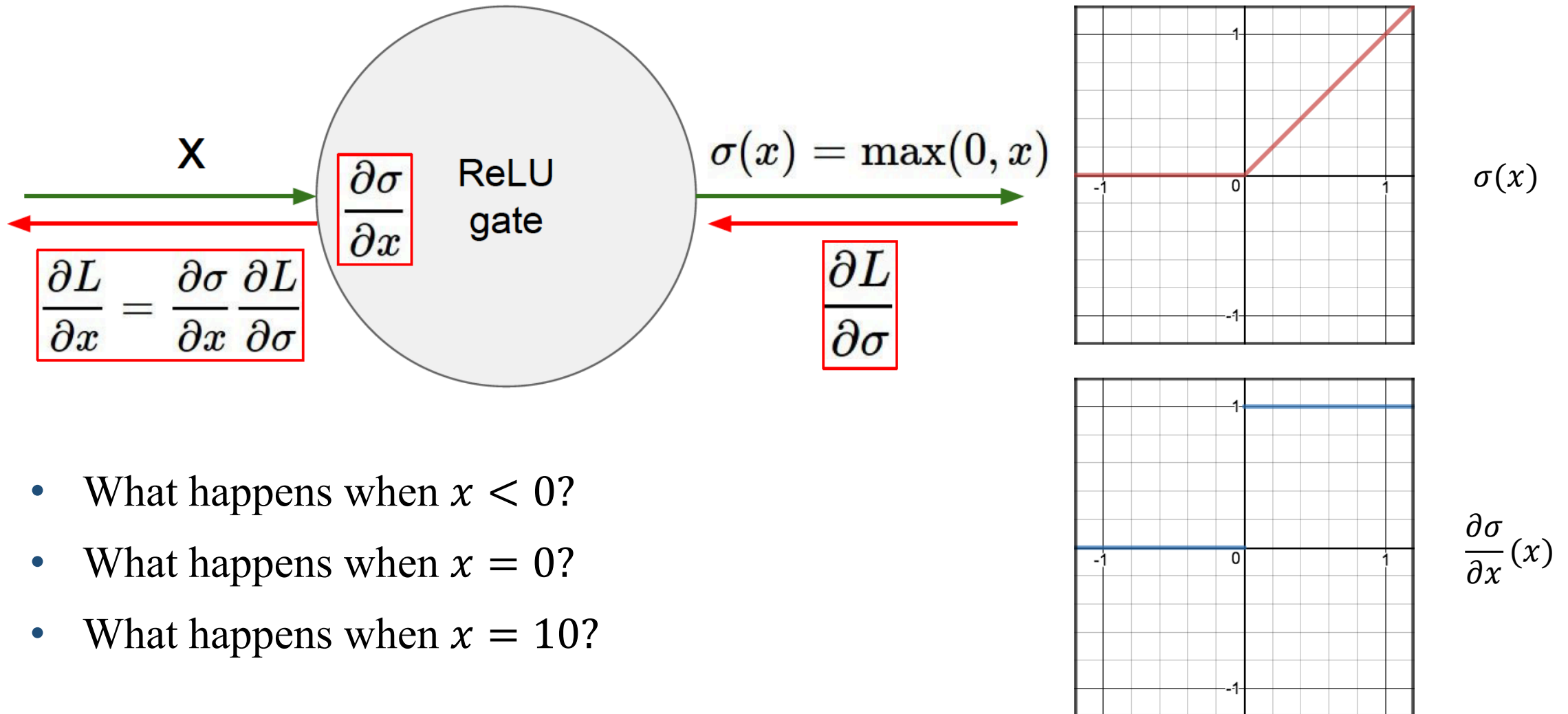
ReLU
(Rectified Linear Unit)

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

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- **Not zero-centered output**

What is the gradient when $x < 0$?

Activation Functions: ReLU

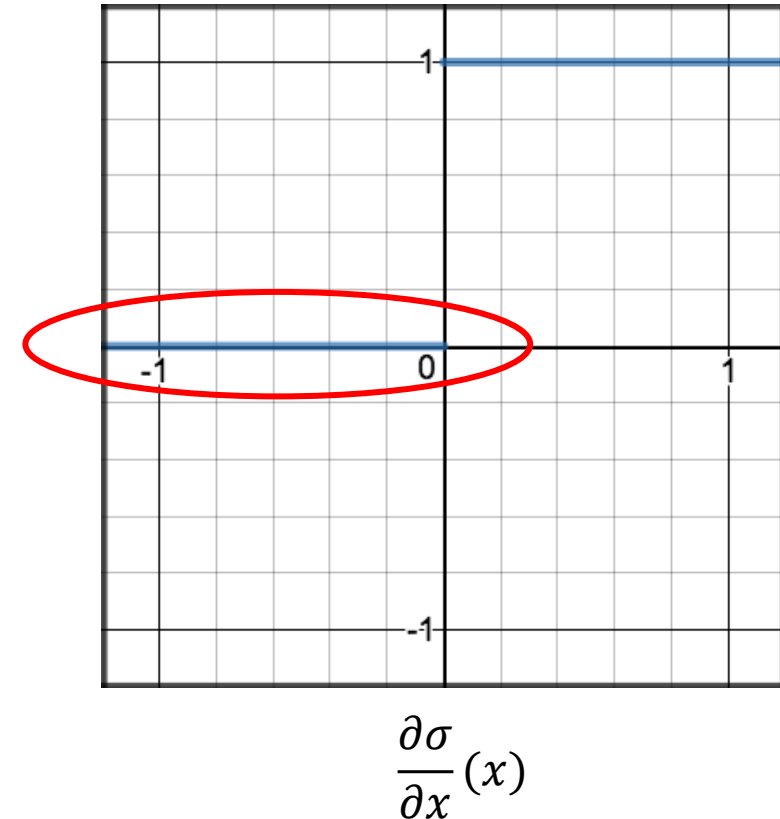


- What happens when $x < 0$?
- What happens when $x = 0$?
- What happens when $x = 10$?

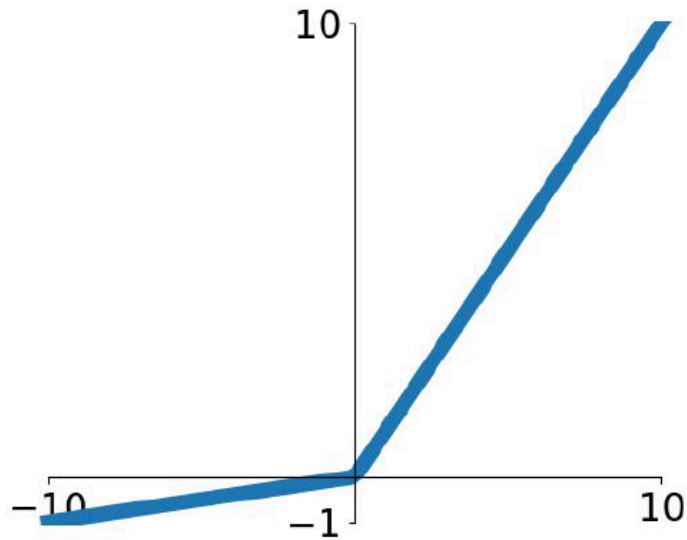
Activation Functions: ReLU

- The dying ReLU Problem:
 - When the input to ReLU is negative, the gradient will be **zero**!
 - Zero gradient will not produce any updates
 - Once you get stuck at zero gradient, you will never get your weights updated again!

⇒ people like to initialize ReLU neurons with slightly positive intercept w_0 (e.g. 0.01)



Activation Functions: Leaky ReLU



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

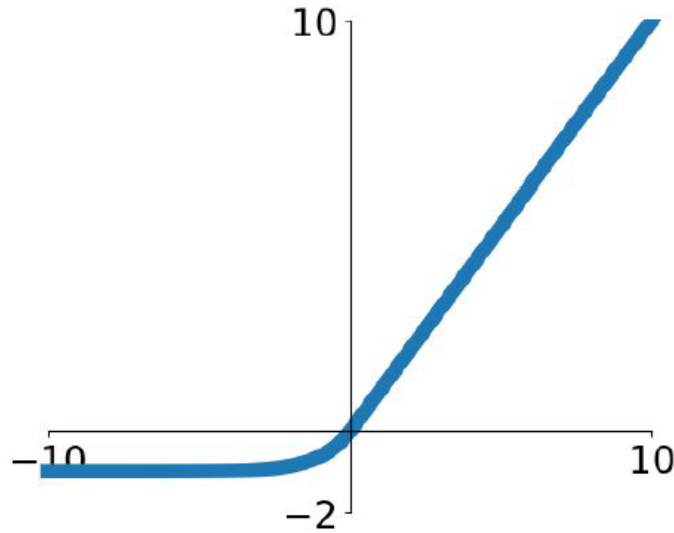
Leaky ReLU

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

backprop into α

By default $\alpha = 0.01$

Activation Functions: ELU



Exponential Linear Units (ELU)

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- All benefits of ReLU
- Closer to zero mean outputs
- Will not “die”
- Saturate on the negative regime
- Computation requires `exp()`