

CS 6316 Machine Learning

Neural Networks

Yangfeng Ji

Information and Language Processing Lab
Department of Computer Science
University of Virginia



1. From Perceptrons to MLPs
2. From Logistic Regression to Neural Networks
3. Expressive Power of Neural Networks
4. Learning Neural Networks
5. Computation Graph

From Perceptrons to MLPs

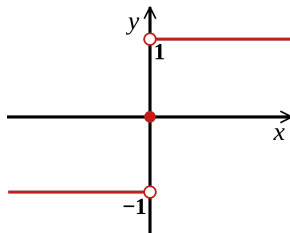
Perceptrons

- ▶ $\mathcal{X} = \mathbb{R}^d$
- ▶ $\mathcal{Y} = \{-1, +1\}$
- ▶ Halfspace hypothesis class

$$\mathcal{H}_{\text{half}} = \{\text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d\} \quad (1)$$

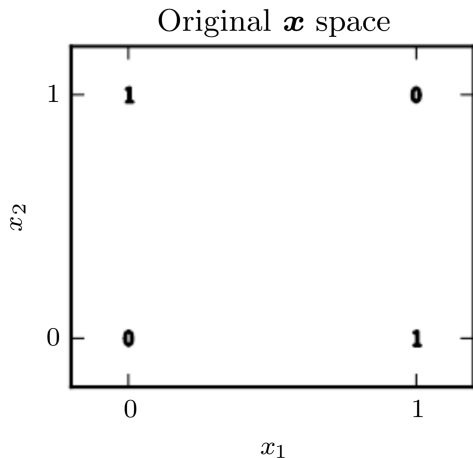
which is an **infinite** hypothesis space.

The sign function $y = \text{sign}(x)$ is defined as



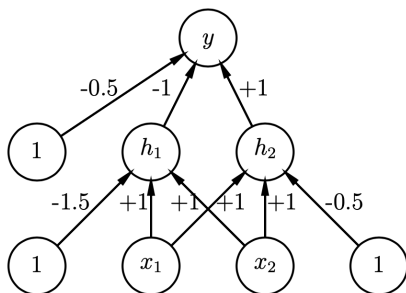
The XOR Problem

$$y = x_1 \oplus x_2 \quad (2)$$



A Multi-Layer Perceptron

The problem can be solved by stacking three perceptrons together, for example,



The new model is called Multi-Layer Perceptron (MLP).

Geometric Interpretation

The previous MLP can be write in the mathematical form as

$$h_1 = \text{sign}(x_1 + x_2 - 1.5) \quad (3)$$

$$h_2 = \text{sign}(x_1 + x_2 - 0.5) \quad (4)$$

$$y = \text{sign}(-h_1 + h_2 - 0.5) \quad (5)$$

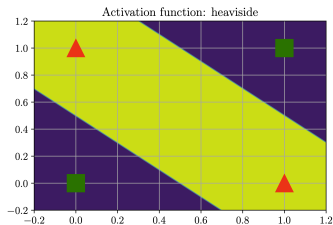
Geometric Interpretation

The previous MLP can be written in the mathematical form as

$$h_1 = \text{sign}(x_1 + x_2 - 1.5) \quad (3)$$

$$h_2 = \text{sign}(x_1 + x_2 - 0.5) \quad (4)$$

$$y = \text{sign}(-h_1 + h_2 - 0.5) \quad (5)$$



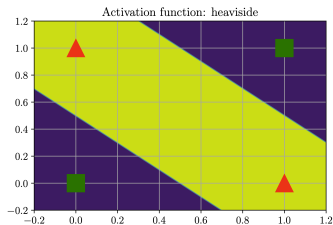
Geometric Interpretation

The previous MLP can be written in the mathematical form as

$$h_1 = \text{sign}(x_1 + x_2 - 1.5) \quad (3)$$

$$h_2 = \text{sign}(x_1 + x_2 - 0.5) \quad (4)$$

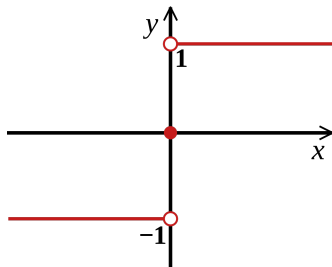
$$y = \text{sign}(-h_1 + h_2 - 0.5) \quad (5)$$



- ▶ Each h_i defines a classifier by dividing the input space into two half-spaces
- ▶ Equation 3 forms a non-linear classifier by combining two linear classifiers together

What about Learning?

- ▶ Although the previous classifier is simple and intuitive, learning the parameters are not easy, because function $\text{sign}(\cdot)$ is non-differentiable!



- ▶ *Solution:* replace $\text{sign}(\cdot)$ function with the Sigmoid function $\sigma(\cdot)$
 - ▶ For example, $h_1 = \sigma(w_1x_1 + w_2x_2)$
 - ▶ In other words, transform each perceptron classifier to a logistic regression classifier

From Logistic Regression to Neural Networks

- ▶ An unified form for $y \in \{-1, +1\}$

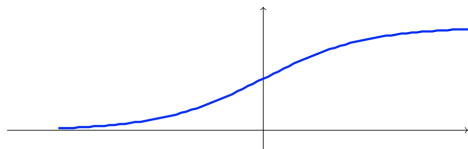
$$p(Y = +1 | \mathbf{x}) = \frac{1}{1 + \exp(-\langle \mathbf{w}, \mathbf{x} \rangle)} \quad (6)$$

- ▶ An unified form for $y \in \{-1, +1\}$

$$p(Y = +1 | \mathbf{x}) = \frac{1}{1 + \exp(-\langle \mathbf{w}, \mathbf{x} \rangle)} \quad (6)$$

- ▶ The sigmoid function $\sigma(a)$ with $a \in \mathbb{R}$

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \quad (7)$$



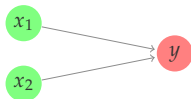
Graphical Representation

- ▶ A specific example of LR

$$p(Y = 1 | \mathbf{x}) = \sigma\left(\sum_{j=1}^2 w_j \mathbf{x}_{.,j}\right) \quad (8)$$

- ▶ The graphical representation of this LR model is

Input layer Output layer



From LR to Neural Networks

Build upon logistic regression, a simple neural network can be constructed as

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (9)$$

$$P(y = 1 | \mathbf{x}) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (10)$$

- ▶ $\mathbf{x} \in \mathbb{R}^d$: d -dimensional input
- ▶ $y \in \{-1, +1\}$ (binary classification problem)

From LR to Neural Networks

Build upon logistic regression, a simple neural network can be constructed as

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{.,j}\right) \quad k \in [K] \quad (9)$$

$$P(y = 1 | \mathbf{x}) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (10)$$

- ▶ $\mathbf{x} \in \mathbb{R}^d$: d -dimensional input
- ▶ $y \in \{-1, +1\}$ (binary classification problem)
- ▶ $\{w_{k,i}^{(1)}\}$ and $\{w_k^{(o)}\}$ are two sets of the parameters, and
- ▶ K is the number of hidden units, each of them has the same form as a LR.

- ▶ Element-wise formulation

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (11)$$

$$P(y = +1 | \mathbf{x}) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (12)$$

- ▶ Element-wise formulation

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (11)$$

$$P(y = +1 | \mathbf{x}) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (12)$$

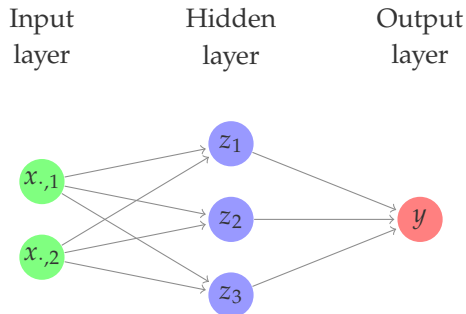
- ▶ Matrix-vector formulation

$$\mathbf{z} = \sigma(\mathbf{W}^{(1)} \mathbf{x}) \quad (13)$$

$$P(y = +1 | \mathbf{x}) = \sigma((\mathbf{w}^{(o)})^\top \mathbf{z}) \quad (14)$$

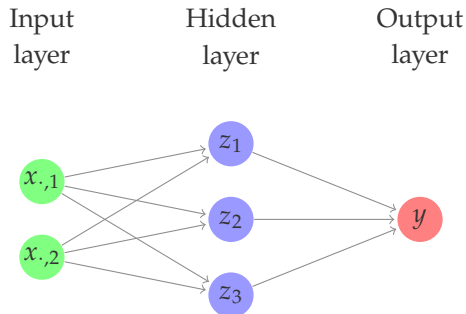
where $\mathbf{W}^{(1)} \in \mathbb{R}^{K \times d}$ and $\mathbf{w}^{(o)} \in \mathbb{R}^K$

Graphical Representation



- ▶ Depth: 2 (two-layer neural network)
- ▶ Width: 3 (the maximal number of units in each layer)

Graphical Representation



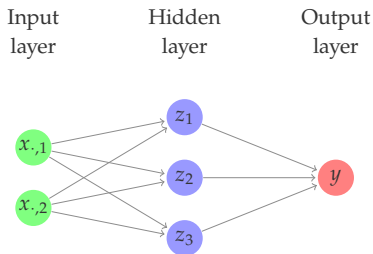
- ▶ Depth: 2 (two-layer neural network)
- ▶ Width: 3 (the maximal number of units in each layer)

Demo for solve the XOR problem

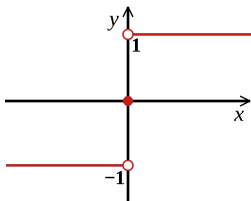
Hypothesis Space

The hypothesis space of neural networks is usually defined by the **architecture** of the network, which includes

- ▶ the nodes in the network,
- ▶ the connections in the network, and
- ▶ the **activation function** (e.g., σ , \tanh)

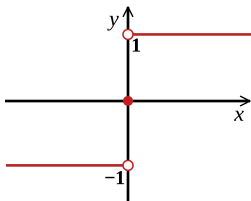


Other Activation Functions

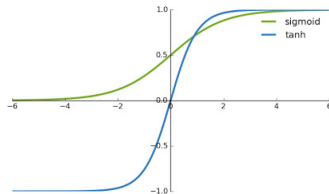


(a) Sign function

Other Activation Functions

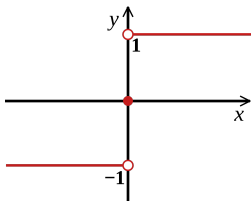


(a) Sign function

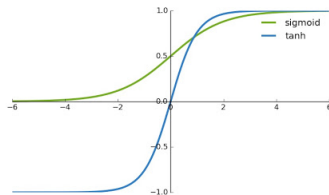


(b) Tanh function

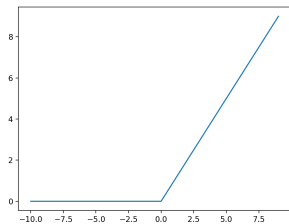
Other Activation Functions



(a) Sign function



(b) Tanh function



(c) ReLU function
[Jarrett et al., 2009]

Expressive Power of Neural Networks

Two-layer NNs with Sign Function

Consider a neural network defined by the following functions

$$z_k = \text{sign}\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (15)$$

$$h(\mathbf{x}) = \text{sign}\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (16)$$

where $\text{sign}(a)$ is the sign function.

Two-layer NNs with Sign Function

Consider a neural network defined by the following functions

$$z_k = \text{sign}\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (15)$$

$$h(\mathbf{x}) = \text{sign}\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (16)$$

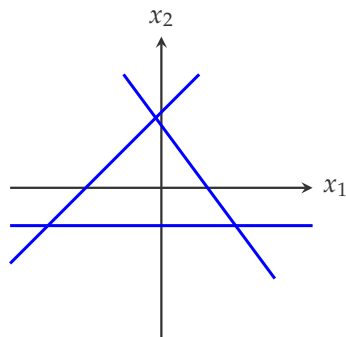
where $\text{sign}(a)$ is the sign function.

$h(\mathbf{x})$ can be rewritten as

$$h(\mathbf{x}) = \text{sign}\left(\sum_{k=1}^K w_k^{(o)} \cdot \text{sign}\left(\sum_{j=1}^d w_{k,i}^{(1)} x_{\cdot,j}\right)\right) \quad (17)$$

Decision Boundary

$h(\mathbf{x})$ is defined by a combination of K linear predictors



Similar conclusion applies to other activation functions. [Demo]

[Shalev-Shwartz and Ben-David, 2014, Page 274]

Universal Approximation Theorem

Restrict the inputs $x_{.,j} \in \{-1, +1\} \forall j \in [d]$ as binary

Universal Approximation Theorem

For every d , there exists a two-layer neural network (Equations 15 – 16), such that this hypothesis space contains all functions from $\{-1, +1\}^d$ to $\{-1, +1\}$

[Shalev-Shwartz and Ben-David, 2014, Section 20.3]

Universal Approximation Theorem

Restrict the inputs $x_{.,j} \in \{-1, +1\} \forall j \in [d]$ as binary

Universal Approximation Theorem

For every d , there exists a two-layer neural network (Equations 15 – 16), such that this hypothesis space contains all functions from $\{-1, +1\}^d$ to $\{-1, +1\}$

- ▶ The minimal size of network that satisfies the theorem is **exponential** in d
- ▶ Similar results hold for σ as the activation function

[Shalev-Shwartz and Ben-David, 2014, Section 20.3]

Learning Neural Networks

Neural Network Predictions

Consider a binary classification problem with $\mathcal{Y} = \{-1, +1\}$,

- ▶ A two-layer neural network gives the following prediction as

$$P(Y = +1 | \mathbf{x}) = \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \quad (18)$$

where $\{\mathbf{w}^{(0)}, \mathbf{W}^{(1)}\}$ are the parameters

Neural Network Predictions

Consider a binary classification problem with $\mathcal{Y} = \{-1, +1\}$,

- ▶ A two-layer neural network gives the following prediction as

$$P(Y = +1 | \mathbf{x}) = \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \quad (18)$$

where $\{\mathbf{w}^{(0)}, \mathbf{W}^{(1)}\}$ are the parameters

- ▶ Assume the ground-truth label is y , let's introduce an **empirical** distribution

$$q(Y = y' | \mathbf{x}) = \delta(y', y) = \begin{cases} 1 & y' = y \\ 0 & y' \neq y \end{cases} \quad (19)$$

Given one data point, The loss function of a neural network is usually defined as the **cross entropy** of the prediction distribution p and the empirical distribution q

$$\begin{aligned} H(q, p) = & -q(Y = +1 | \mathbf{x}) \log p(Y = +1 | \mathbf{x}) \\ & -q(Y = -1 | \mathbf{x}) \log p(Y = -1 | \mathbf{x}) \end{aligned} \quad (20)$$

Given one data point, The loss function of a neural network is usually defined as the **cross entropy** of the prediction distribution p and the empirical distribution q

$$\begin{aligned} H(q, p) = & -q(Y = +1 | \mathbf{x}) \log p(Y = +1 | \mathbf{x}) \\ & -q(Y = -1 | \mathbf{x}) \log p(Y = -1 | \mathbf{x}) \end{aligned} \quad (20)$$

Since q is defined with a Delta function, Depending on y , we have

$$H(q, p) = \begin{cases} -\log p(Y = +1 | \mathbf{x}) & Y = +1 \\ -\log p(Y = -1 | \mathbf{x}) & Y = -1 \end{cases} \quad (21)$$

Given one data point, The loss function of a neural network is usually defined as the **cross entropy** of the prediction distribution p and the empirical distribution q

$$\begin{aligned} H(q, p) = & -q(Y = +1 | \mathbf{x}) \log p(Y = +1 | \mathbf{x}) \\ & -q(Y = -1 | \mathbf{x}) \log p(Y = -1 | \mathbf{x}) \end{aligned} \quad (20)$$

Since q is defined with a Delta function, Depending on y , we have

$$H(q, p) = \begin{cases} -\log p(Y = +1 | \mathbf{x}) & Y = +1 \\ -\log p(Y = -1 | \mathbf{x}) & Y = -1 \end{cases} \quad (21)$$

It is equivalent to the negative log-likelihood (NLL) function used in learning LR.

- ▶ Given a set of training example $S = \{(x_i, y_i)\}_{i=1}^m$, the loss function is defined as

$$L(\theta) = - \sum_{i=1}^m \log p(y_i | x_i) \quad (22)$$

where θ indicates all the parameters in a network.

- ▶ Given a set of training example $S = \{(x_i, y_i)\}_{i=1}^m$, the loss function is defined as

$$L(\theta) = - \sum_{i=1}^m \log p(y_i | x_i) \quad (22)$$

where θ indicates all the parameters in a network.

- ▶ For example, $\theta = \{w^{(o)}, \mathbf{W}^{(1)}\}$, for the previously defined two-layer neural network

- ▶ Given a set of training example $S = \{(x_i, y_i)\}_{i=1}^m$, the loss function is defined as

$$L(\theta) = - \sum_{i=1}^m \log p(y_i | x_i) \quad (22)$$

where θ indicates all the parameters in a network.

- ▶ For example, $\theta = \{w^{(0)}, \mathbf{W}^{(1)}\}$, for the previously defined two-layer neural network
- ▶ Just like learning a LR, we can use **gradient-based** learning algorithm

Gradient-based Learning

A simple scratch of gradient-based learning¹

1. Compute the gradient of θ , $\frac{\partial L(\theta)}{\partial \theta}$

¹More detail will be discussed in the next lecture

A simple scratch of gradient-based learning¹

1. Compute the gradient of θ , $\frac{\partial L(\theta)}{\partial \theta}$
2. Update the parameter with the gradient

$$\theta^{(\text{new})} \leftarrow \theta^{(\text{old})} - \eta \cdot \left. \frac{\partial L(\theta)}{\partial \theta} \right|_{\theta = \theta^{(\text{old})}} \quad (23)$$

where η is the learning rate

¹More detail will be discussed in the next lecture

A simple scratch of gradient-based learning¹

1. Compute the gradient of θ , $\frac{\partial L(\theta)}{\partial \theta}$
2. Update the parameter with the gradient

$$\theta^{(\text{new})} \leftarrow \theta^{(\text{old})} - \eta \cdot \left. \frac{\partial L(\theta)}{\partial \theta} \right|_{\theta=\theta^{(\text{old})}} \quad (23)$$

where η is the learning rate

3. Go back step 1 until it converges

¹More detail will be discussed in the next lecture

Consider the two-layer neural network with one training example (\mathbf{x}, y) , to further simplify the computation, we assume $y = +1$

$$\log p(y | \mathbf{x}) = \log \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \quad (24)$$

Gradient Computation

Consider the two-layer neural network with one training example (\mathbf{x}, y) , to further simplify the computation, we assume $y = +1$

$$\log p(y | \mathbf{x}) = \log \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \quad (24)$$

The gradient with respect to $\mathbf{w}^{(0)}$ is

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \mathbf{w}^{(0)}} = - \frac{\partial \log \sigma(\cdot)}{\partial \sigma(\cdot)} \cdot \frac{\partial \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right)}{\partial (\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x})} \cdot \frac{\partial (\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x})}{\partial \mathbf{w}^{(0)}} \quad (25)$$

Gradient Computation

Consider the two-layer neural network with one training example (\mathbf{x}, y) , to further simplify the computation, we assume $y = +1$

$$\log p(y | \mathbf{x}) = \log \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \quad (24)$$

The gradient with respect to $\mathbf{w}^{(0)}$ is

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \mathbf{w}^{(0)}} &= - \frac{\partial \log \sigma(\cdot)}{\partial \sigma(\cdot)} \cdot \frac{\partial \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right)}{\partial (\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x})} \cdot \frac{\partial (\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x})}{\partial \mathbf{w}^{(0)}} \\ &= - \left\{ 1 - \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \right\} \cdot \sigma(\mathbf{W}^{(1)} \mathbf{x}) \end{aligned} \quad (25)$$

which is in the similar form as the LR updating equation.

The gradient with respect to $W^{(1)}$ is

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \mathbf{W}^{(1)}} &= \frac{\partial \log \sigma(\cdot)}{\partial \sigma(\cdot)} \cdot \frac{\partial \sigma\left((\mathbf{w}^{(o)})^\top \sigma(\mathbf{W}^{(1)}\mathbf{x})\right)}{\partial (\mathbf{w}^{(o)})^\top \sigma(\mathbf{W}^{(1)}\mathbf{x})} \\ &\quad \cdot \frac{\partial (\mathbf{w}^{(o)})^\top \sigma(\mathbf{W}^{(1)}\mathbf{x})}{\partial \sigma(\mathbf{W}^{(1)}\mathbf{x})} \cdot \frac{\partial \sigma(\mathbf{W}^{(1)}\mathbf{x})}{\partial \mathbf{W}^{(1)}\mathbf{x}} \cdot \frac{\partial \mathbf{W}^{(1)}\mathbf{x}}{\partial \mathbf{W}^{(1)}} \end{aligned} \quad (26)$$

The gradient with respect to $W^{(1)}$ is

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \mathbf{W}^{(1)}} &= \frac{\partial \log \sigma(\cdot)}{\partial \sigma(\cdot)} \cdot \frac{\partial \sigma((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}))}{\partial (\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x})} \\ &\quad \cdot \frac{\partial (\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x})}{\partial \sigma(\mathbf{W}^{(1)} \mathbf{x})} \cdot \frac{\partial \sigma(\mathbf{W}^{(1)} \mathbf{x})}{\partial \mathbf{W}^{(1)} \mathbf{x}} \cdot \frac{\partial \mathbf{W}^{(1)} \mathbf{x}}{\partial \mathbf{W}^{(1)}} \end{aligned} \quad (26)$$

- ▶ Both of them are the applications of the chain rule in calculus plus some derivatives of basic functions
- ▶ In the literature of neural networks, it is called the **back-propagation** algorithm [Rumelhart et al., 1986]

Computation Graph

Consider the example of a two-layer neural network

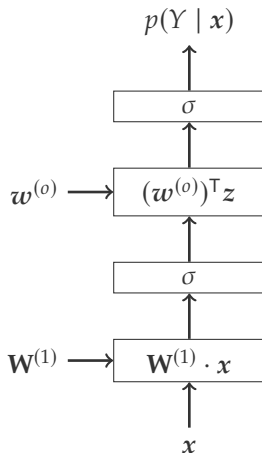
$$P(Y = +1 | \mathbf{x}) = \sigma \left((\mathbf{w}^{(0)})^\top \sigma(\mathbf{W}^{(1)} \mathbf{x}) \right) \quad (27)$$

A neural network is a composition of some basic functions and operations. For example

- ▶ $\sigma(\cdot)$
- ▶ matrix transpose $(\mathbf{w}^{(0)})^\top$
- ▶ matrix-vector multiplication $\mathbf{W}^{(1)} \mathbf{x}$

Forward Graph

The computation graph of the two-layer neural network²



²For simplicity, the transpose operation is ignored from the graph

Backward Operations

Similarly, the gradient of neural network parameters are computed with a series of backward operations associated with the derivative of some basic function. For example

$$\blacktriangleright \frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$

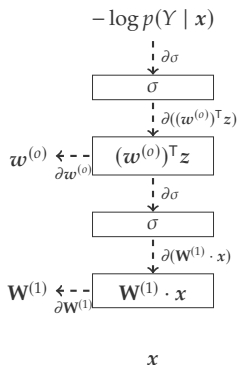
$$\blacktriangleright \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\blacktriangleright \frac{\partial \log(x)}{\partial x} = \frac{1}{x}$$

$$\blacktriangleright \frac{\partial \mathbf{W}\mathbf{x}}{\partial \mathbf{x}} = \begin{bmatrix} \mathbf{x}^\top \\ \vdots \\ \mathbf{x}^\top \end{bmatrix}$$

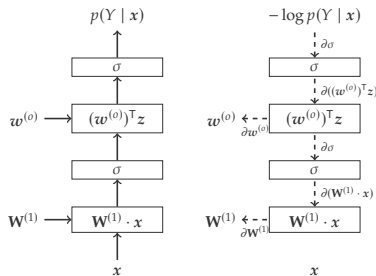
Backward Graph

With the chain rule, gradient of the loss function with respect to any parameter can be computed backward step-by-step along the path



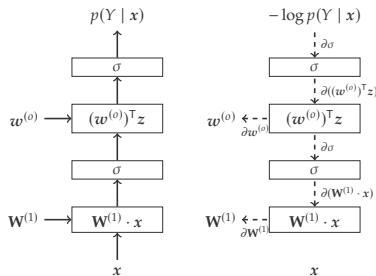
Computation Graph

Perform the forward/backward step with a graph of basic operations
(e.g., PyTorch, Tensorflow)



Computation Graph

Perform the forward/backward step with a graph of basic operations (e.g., PyTorch, Tensorflow)



- ▶ Modular implementation: implement each module with its forward/backward operations together
- ▶ Automatic differentiation: automatically run with the backward step

What is Deep Learning?

Definition

Deep Learning is building a system by assembling **parameterized modules** into a (possibly dynamic) **computation graph**, and training it to perform a task by optimizing the parameters using a **gradient-based method**.

[LeCun, 2020, AAAI 2020 Keynote]



Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2009).
What is the best multi-stage architecture for object recognition?
In *Proceedings of the 12th International Conference on Computer Vision*, pages 2146–2153. IEEE.



LeCun, Y. (2020).
Self-supervised learning.



Pascanu, R., Mikolov, T., and Bengio, Y. (2013).
On the difficulty of training recurrent neural networks.
In Dasgupta, S. and McAllester, D., editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1310–1318, Atlanta, Georgia, USA. PMLR.



Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986).
Learning representations by back-propagating errors.
Nature, 323(6088):533–536.



Shalev-Shwartz, S. and Ben-David, S. (2014).
Understanding machine learning: From theory to algorithms.
Cambridge university press.