

# Introduction to Deep Learning

Jiong Zhang

Apr 2020

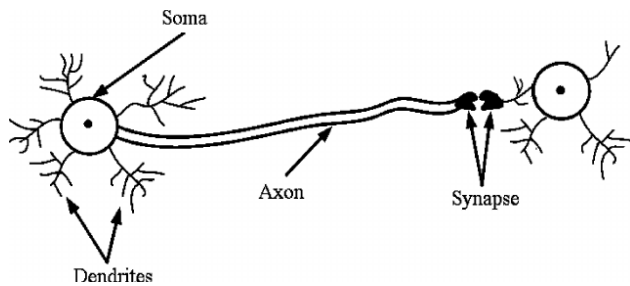
## 1 Overview

**Example 1.** Computer Vision (CV)

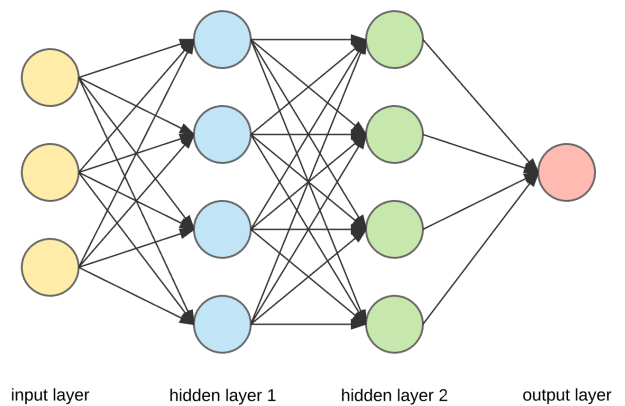
- object detection
- face recognition
- self-driving cars

**Example 2.** Natural Language Processing (NLP)

- machine translation
- article summarization
- search suggestion



Biological Neural Networks



Artificial Neural Networks

Model	Neurons	Synapses
Human	8.6e10	1.5e14
Cat	7.6e8	1.0e13
Bee	1.0e6	1.0e9
ResNet101	$\sim 1.0e5$	4.4e7
BERT(large) <sup>1</sup>	$\sim 1.0e6$	1.1e8

Table 1: Size comparison of Biological Neural Networks and Artificial Neural Networks

## 2 Feed-forward Neural Networks

### 2.1 Matrix Differentiation

Preliminaries for FNN back propagation derivation. Two major layouts for matrix differentiation:

- numerator layout:

$$\frac{\partial y}{\partial \vec{x}} = \left[ \frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \dots \quad \frac{\partial y}{\partial x_n} \right] \quad \frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathcal{R}^{m \times n} \quad (1)$$

- denominator layout

$$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \dots \\ \frac{\partial y}{\partial x_n} \end{bmatrix} \quad \frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathcal{R}^{n \times m} \quad (2)$$

To simplify the notation, we adopt the numerator layout.

#### 2.1.1 Basic Definitions

- scalar/vector:  $y \in \mathcal{R}$ ,  $\mathbf{x} \in \mathcal{R}^n$

$$\frac{\partial y}{\partial \vec{x}} = \left[ \frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \dots \quad \frac{\partial y}{\partial x_n} \right]; \left[ \frac{\partial y}{\partial \vec{x}} \right]^\top \in \mathcal{R}^n \quad (3)$$

- vector/scalar:  $\mathbf{y} \in \mathcal{R}^m$ ,  $x \in \mathcal{R}$

$$\frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \dots \\ \frac{\partial y_m}{\partial x} \end{bmatrix} \in \mathcal{R}^m \quad (4)$$

- vector/vector:  $\mathbf{y} \in \mathcal{R}^m$ ,  $\mathbf{x} \in \mathcal{R}^n$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathcal{R}^{m \times n} \quad (5)$$

- scalar/matrix:  $y \in \mathcal{R}$ ,  $\mathbf{X} \in \mathcal{R}^{m \times n}$

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{1,1}} & \dots & \frac{\partial y}{\partial X_{m,1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial X_{1,n}} & \dots & \frac{\partial y}{\partial X_{m,n}} \end{bmatrix} \in \mathcal{R}^{n \times m} \quad (6)$$

- $a$  order tensor/ $b$  order tensor:  $\mathbf{Y} \in \mathcal{R}^{m_1 \times m_2 \times \dots \times m_a}$ ,  $\mathbf{X} \in \mathcal{R}^{n_1 \times n_2 \times \dots \times n_b}$

$$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}} \in \mathcal{R}^{m_1 \times \dots \times m_a \times n_b \times \dots \times n_1} \quad (7)$$

### 2.1.2 Derivative Formulas

**Case 1.**  $\vec{y} = \mathbf{A}\vec{x}$ ,  $\vec{x} \in \mathcal{R}^n$ ,  $\mathbf{A} \in \mathcal{R}^{m \times n}$

$$\left[ \frac{\partial \vec{y}}{\partial \vec{x}} \right]_{i,j} = \frac{\partial y_i}{\partial x_j} = \frac{\partial \sum_{t=1}^n A_{i,t} x_t}{\partial x_j} = A_{i,j} \quad (8)$$

Therefore

$$\frac{\partial \mathbf{A}\vec{x}}{\partial \vec{x}} = \mathbf{A} \quad (9)$$

**Case 2.**  $\vec{y} = \mathbf{X}\vec{a}$ ,  $\mathbf{X} \in \mathcal{R}^{m \times n}$ ,  $\vec{a} \in \mathcal{R}^n$

$$\left[ \frac{\partial \vec{y}}{\partial \mathbf{X}} \right]_{i,j,k} = \frac{\partial y_i}{\partial X_{k,j}} = \frac{\partial \sum_{t=1}^n X_{i,t} a_t}{\partial X_{k,j}} = a_j \mathbb{I}(i=k) \quad (10)$$

Therefore

$$\frac{\partial \vec{y}}{\partial \mathbf{X}} = \rho(\vec{a}, m) \in \mathcal{R}^{m \times n \times m} \quad (11)$$

$$\rho(\vec{a}, m)_{i,j,k} \equiv a_j \mathbb{I}(i=k) \quad (12)$$

**Case 3.** Chain Rule:  $\vec{y} \in \mathcal{R}^m$  is a function of  $\vec{z} \in \mathcal{R}^r$  and  $\vec{z}$  is a function of  $\vec{x} \in \mathcal{R}^n$ .

$$\left[ \frac{\partial \vec{y}}{\partial \vec{x}} \right]_{i,j} = \frac{\partial y_i}{\partial x_j} \quad (13)$$

By chain rule

$$\frac{\partial y_i}{\partial x_j} = \sum_{t=1}^r \frac{\partial y_i}{\partial z_t} \frac{\partial z_t}{\partial x_j} = \frac{\partial y_i}{\partial \vec{z}} \frac{\partial \vec{z}}{\partial x_j} = \left[ \frac{\partial \vec{y}}{\partial \vec{z}} \right]_{i,:} \left[ \frac{\partial \vec{z}}{\partial \vec{x}} \right]_{:,j} \quad (14)$$

Therefore

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \frac{\partial \vec{y}}{\partial \vec{z}} \frac{\partial \vec{z}}{\partial \vec{x}} \quad (15)$$

Chain rule of matrix differentiation (numerator layout) "toward right". For denominator layout it is "toward left".

## 2.2 Multi-layer Perceptrons: neural networks

Previously we learned about perceptrons, however, in most real applications, a linear combination of input features is not enough to make accurate predictions. We need to capture interactions among input features (feature extraction).

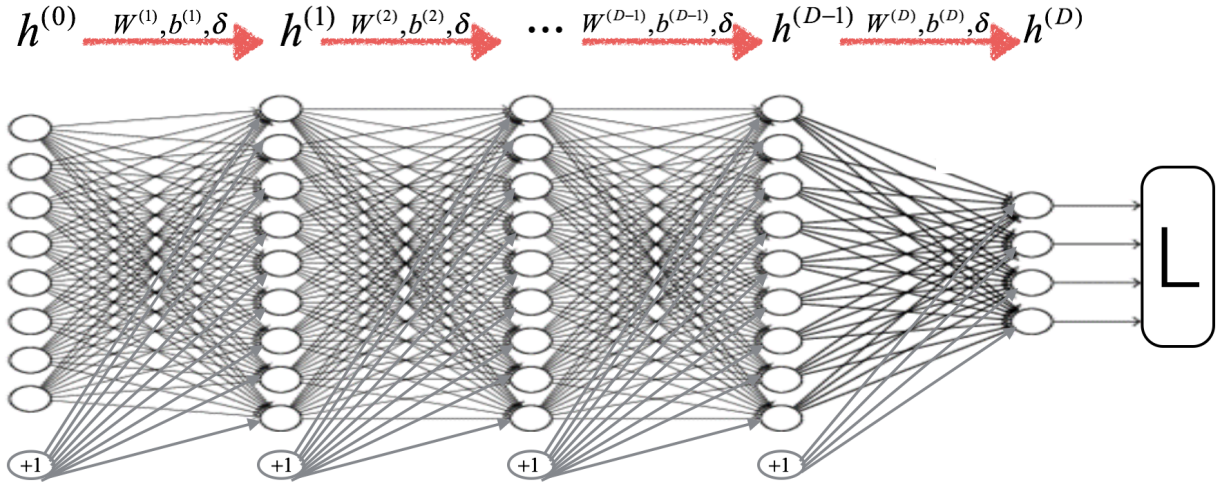


Figure 1: A Multi-layer Perceptron with depth  $D$

Consider regression problem with input data  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^M$  and  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^M$ , where  $\mathbf{x}_i \in \mathcal{R}^{n_0}$  and  $\mathbf{y}_i \in \mathcal{R}^{n_D}$ . A Multi-layer Perceptron (MLP) with depth  $D$  is defined as:

$$\mathbf{h}^{(d)} = \delta(\mathbf{W}^{(d)}\mathbf{h}^{(d-1)} + \mathbf{b}^{(d)}), \text{ for } d = 1, 2, \dots, D \quad (16)$$

$$\mathbf{h}^{(0)} = \mathbf{x} \quad (17)$$

where  $\mathbf{h}^{(d)} \in \mathcal{R}^{n_d}$  are the hidden vectors,  $\mathbf{W}^{(d)} \in \mathcal{R}^{n_d \times n_{d-1}}$  and  $\mathbf{b}^{(d)} \in \mathcal{R}^{n_d}$  are trainable parameters. The loss function is defined on the final outputs  $\mathbf{h}^{(D)}$ :

$$L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Y}) = \frac{1}{M} \sum_{i=1}^M \hat{l}(\boldsymbol{\theta}; \mathbf{x}_i, \mathbf{y}_i) = \frac{1}{M} \sum_{i=1}^M l(\mathbf{h}^{(D)}(\boldsymbol{\theta}, \mathbf{x}_i), \mathbf{y}_i) \quad (18)$$

where  $\boldsymbol{\theta} = [\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(D)}, \mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \dots, \mathbf{b}^{(D)}]$ .

### Model Training:SGD

- Initialize  $\mathbf{W}_0^{(d)}, \mathbf{b}_0^{(d)}$  for  $d = 1, 2, \dots, D$ , learning rate  $\eta \in \mathcal{R}$
- For  $t = 1, 2, \dots, T$ :
  - Sample  $i \sim \text{Uniform}(1, 2, \dots, M)$
  - For  $d = D, D - 1, \dots, 1$ :
    - \*  $\mathbf{W}_{t+1}^{(d)} \leftarrow \mathbf{W}_t^{(d)} - \eta \left[ \frac{\partial \hat{l}}{\partial \mathbf{W}^{(d)}}(\boldsymbol{\theta}_t; \mathbf{x}_i, \mathbf{y}_i) \right]^\top$
    - \*  $\mathbf{b}_{t+1}^{(d)} \leftarrow \mathbf{b}_t^{(d)} - \eta \left[ \frac{\partial \hat{l}}{\partial \mathbf{b}^{(d)}}(\boldsymbol{\theta}_t; \mathbf{x}_i, \mathbf{y}_i) \right]^\top$

**Forward Propagation.** Note that each hidden vector  $\mathbf{h}^{(d)}$  is dependent on its previous layer  $\mathbf{h}^{(d-1)}$ . Therefore, to compute the model output  $\mathbf{h}^{(D)}$ , all the hidden vectors needs to be evaluated in a forwarding sequence.

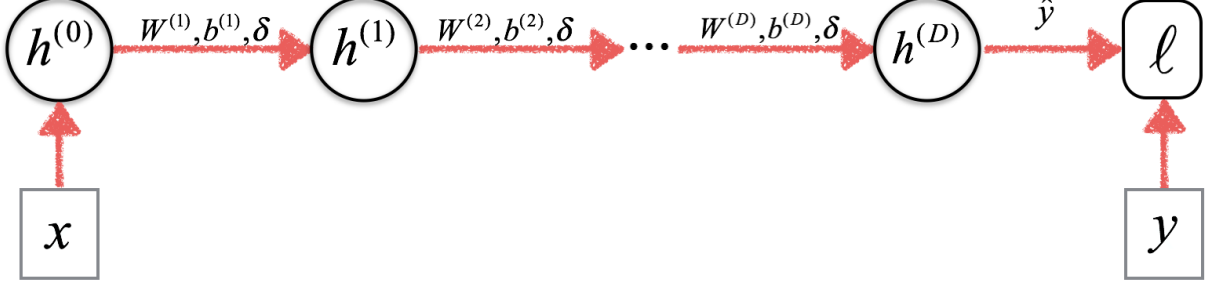


Figure 2: Forward propagation

### 2.3 Back Propagation.

To train parameters  $\mathbf{W}^{(d)}$ , we need to evaluate  $\frac{\partial l}{\partial \mathbf{w}^{(d)}}(\boldsymbol{\theta}_t; \mathbf{x}_i, \mathbf{y}_i)$ . The explicit formulation for the MLP can be written as:

$$l_i = \widehat{l}(\boldsymbol{\theta}; \mathbf{x}_i, \mathbf{y}_i) = l(\delta(\mathbf{W}^{(D)}\delta(\mathbf{W}^{(D-1)}\dots\delta(\mathbf{W}^{(1)}\mathbf{x}_i + \mathbf{b}^{(1)}) + \dots + \mathbf{b}^{(D)}), \mathbf{y}_i). \quad (19)$$

We suppress subscript  $i$  for simplicity. By definition:

$$\frac{\partial l}{\partial \mathbf{h}^{(D)}} = \left[ \frac{\partial l}{\partial h_1^{(D)}} \quad \frac{\partial l}{\partial h_2^{(D)}} \quad \dots \quad \frac{\partial l}{\partial h_{n_D}^{(D)}} \right] \quad (20)$$

Denoting  $\mathbf{z}^{(d)} = \mathbf{W}^{(d)}\mathbf{h}^{(d-1)} + \mathbf{b}^{(d)}$  The derivatives between adjacent layers can be written as

$$\frac{\partial \mathbf{h}^{(d)}}{\partial \mathbf{h}^{(d-1)}} = \frac{\partial \mathbf{h}^{(d)}}{\partial \mathbf{z}^{(d)}} \frac{\partial \mathbf{z}^{(d)}}{\partial \mathbf{h}^{(d-1)}} \quad (21)$$

where

$$\frac{\partial \mathbf{h}^{(d)}}{\partial \mathbf{z}^{(d)}} \equiv \text{diag}(\delta^{(d)}) = \begin{bmatrix} \delta'(z_1^{(d)}) & & \\ & \ddots & \\ & & \delta'(z_{n_d}^{(d)}) \end{bmatrix} \quad (22)$$

$$= \begin{bmatrix} \delta'(\mathbf{W}_{:,1}^{(d)} \cdot \mathbf{h}^{(d-1)} + b_1^{(d)}) & & \\ & \ddots & \\ & & \delta'(\mathbf{W}_{:,n_d}^{(d)} \cdot \mathbf{h}^{(d-1)} + b_{n_d}^{(d)}) \end{bmatrix} \quad (23)$$

And

$$\frac{\partial \mathbf{z}^{(d)}}{\partial \mathbf{h}^{(d-1)}} = \mathbf{W}^{(d)} \in \mathcal{R}^{n_d \times n_{d-1}} \quad (24)$$

Therefore, by chain rule:

$$\frac{\partial l}{\partial \mathbf{W}^{(d)}} = \frac{\partial l}{\partial \mathbf{h}^{(d)}} \frac{\partial \mathbf{h}^{(d)}}{\partial \mathbf{W}^{(d)}} \quad (25)$$

$$= \frac{\partial l}{\partial \mathbf{h}^{(D)}} \frac{\partial \mathbf{h}^{(D)}}{\partial \mathbf{h}^{(D-1)}} \dots \frac{\partial \mathbf{h}^{(d+1)}}{\partial \mathbf{h}^{(d)}} \frac{\partial \mathbf{h}^{(d)}}{\partial \mathbf{W}^{(d)}} \quad (26)$$

Note that  $\frac{\partial \mathbf{h}^{(d)}}{\partial \mathbf{W}^{(d)}} = \text{diag}(\delta^{(d)})\rho(\mathbf{h}^{(d-1)}, n_d) \in \mathcal{R}^{n_d \times n_{d-1} \times n_d}$  is a third order tensor. The gradient w.r.t.  $\mathbf{W}^{(d)}$  can be written as:

$$\frac{\partial l}{\partial \mathbf{W}^{(d)}} = \frac{\partial l}{\partial \mathbf{h}^{(D)}} \text{diag}(\delta^{(D)})\mathbf{W}^{(D)} \dots \text{diag}(\delta^{(d+1)})\mathbf{W}^{(d+1)} \text{diag}(\delta^{(d)})\rho(\mathbf{h}^{(d-1)}, n_d) \quad (27)$$

Therefore, the gradients can be evaluated as:

$$\frac{\partial l}{\partial \mathbf{W}^{(d)}} = \frac{\partial l}{\partial \mathbf{h}^{(D)}} \left[ \prod_{i=D}^{d+1} \text{diag}(\dot{\delta}^{(i)}) \mathbf{W}^{(i)} \right] \text{diag}(\dot{\delta}^{(d)}) \rho(\mathbf{h}^{(d-1)}, n_d) \in \mathcal{R}^{n_{d-1} \times n_d} \quad (28)$$

$$\frac{\partial l}{\partial \mathbf{b}^{(d)}} = \frac{\partial l}{\partial \mathbf{h}^{(D)}} \left[ \prod_{i=D}^{d+1} \text{diag}(\dot{\delta}^{(i)}) \mathbf{W}^{(i)} \right] \text{diag}(\dot{\delta}^{(d)}) \in \mathcal{R}^{1 \times n_d} \quad (29)$$

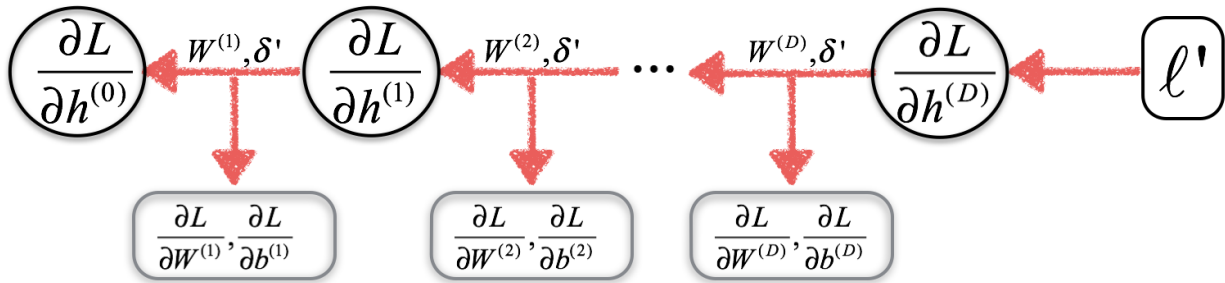


Figure 3: Back propagation

## 2.4 Activation Functions

Activation function	Equation	Example	1D Graph	Derivative
Unit step (Heaviside)	$\phi(z) = \begin{cases} 0, & z < 0, \\ 0.5, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant		
Sign (Signum)	$\phi(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant		
Linear	$\phi(z) = z$	Adaline, linear regression		
Piece-wise linear	$\phi(z) = \begin{cases} 1, & z \geq \frac{1}{2}, \\ z + \frac{1}{2}, & -\frac{1}{2} < z < \frac{1}{2}, \\ 0, & z \leq -\frac{1}{2}, \end{cases}$	Support vector machine		
Logistic (sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Logistic regression, Multi-layer NN		
Hyperbolic tangent	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	Multi-layer Neural Networks		
Rectifier, ReLU (Rectified Linear Unit)	$\phi(z) = \max(0, z)$	Multi-layer Neural Networks		
Rectifier, softplus	$\phi(z) = \ln(1 + e^z)$	Multi-layer Neural Networks		

Copyright © Sebastian Raschka 2016  
(<http://sebastianraschka.com>)

Figure 4: Common activation functions and their derivatives.

- Absolute value of input is large for activations  $\Rightarrow$  Saturated activation functions ( $\delta_i \ll 1$ )  $\Rightarrow$  vanishing gradients
- Solutions: constrain size of each  $h$ , use *ReLU* or *softplus*

### 3 Recurrent Neural Networks

**Question:** What if input data are sequences of different lengths? Eg. translation, voice recognition, stock value prediction.

- MLPs can deal with input data with fixed dimensions.
- For sequence data, recurrent models are usually used.

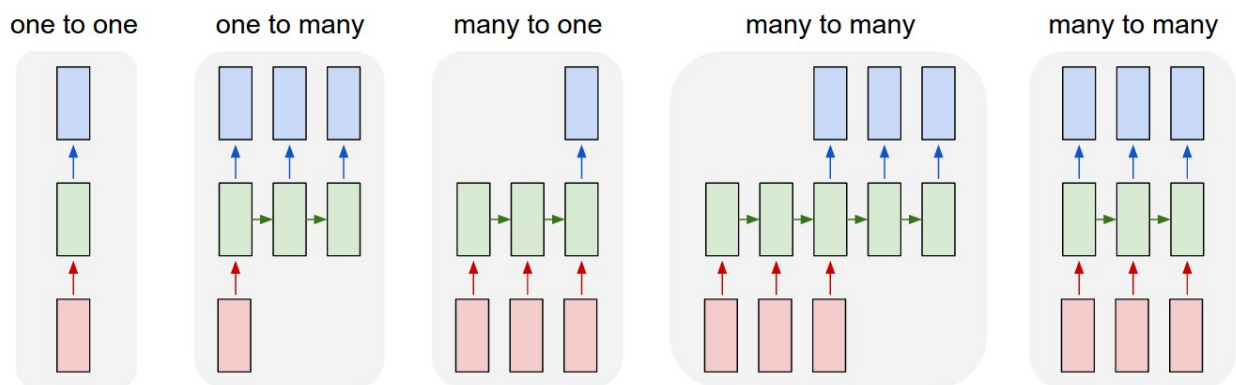


Figure 5: Different use cases of RNNs

**Examples:**

- Neural Machine Translation, chat robot: many to many (case 4)
- Language modelling, time series prediction: many to many (case 5)
- Voice recognition: many to one (case 3)

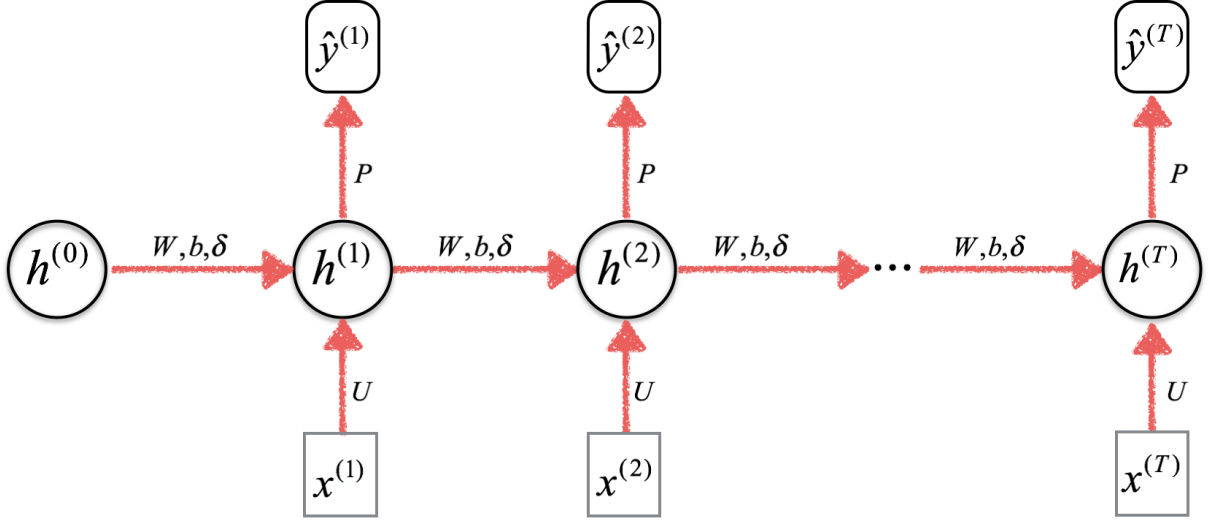


Figure 6: Recurrent Neural Networks: forward propagation

Consider use case 5 with input data  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^M$  and  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^M$ , where

$$\mathbf{x}_i = [\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \dots, \mathbf{x}_i^{(T_i)}], \mathbf{x}_i^{(t)} \in \mathcal{R}^{n_i}$$

$$\mathbf{y}_i = [\mathbf{y}_i^{(1)}, \mathbf{y}_i^{(2)}, \dots, \mathbf{y}_i^{(T_i)}], \mathbf{y}_i^{(t)} \in \mathcal{R}^{n_o}$$

- RNN with activation function  $\delta$ :

$$\mathbf{h}^{(t)} = \delta(\mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t-1)} + \mathbf{b}) \quad (30)$$

$$\hat{\mathbf{y}}^{(t)} = \mathbf{P}\mathbf{h}^{(t)} \quad (31)$$

where  $\mathbf{h}^{(t)} \in \mathcal{R}^{n_h}$  are hidden vectors to encode input information

- Key differences from MLP:
  - Parameters  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{U}, \mathbf{b}, \mathbf{P}\}$  are shared for each cell/time-step
  - Input  $\mathbf{x}_i = [\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \dots, \mathbf{x}_i^{(T_i)}]$  is fed sequentially
  - Output  $\hat{\mathbf{y}} = \{\hat{\mathbf{y}}^{(0)}, \hat{\mathbf{y}}^{(1)}, \dots, \hat{\mathbf{y}}^{(T_i)}\}$  is evaluated at each cell/time-step
- Loss is measured as:

$$L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Y}) = \frac{1}{M} \sum_{i=1}^M \hat{l}(\boldsymbol{\theta}, \mathbf{x}_i, \mathbf{y}_i) \quad (32)$$

$$= \frac{1}{M} \sum_{i=1}^M \sum_{j=1}^{T_i} l(\hat{\mathbf{y}}_i^{(j)}, \mathbf{y}_i^{(j)}) \quad (33)$$



### 3.1 Back propagation

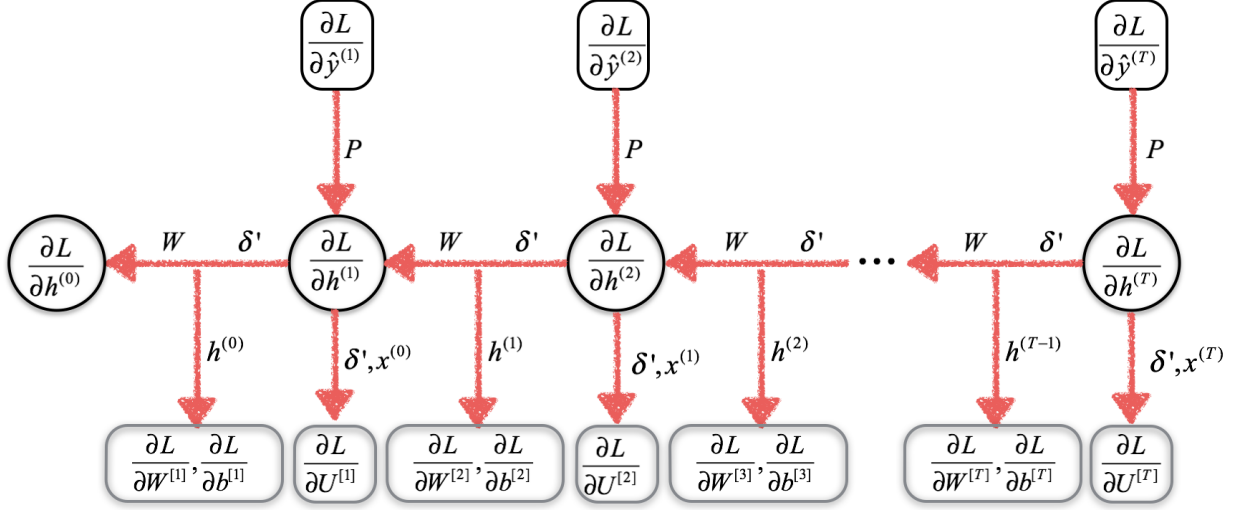


Figure 7: Recurrent Neural Networks: back propagation

Suppressing  $i$ , let  $l_j \equiv l(\hat{\mathbf{y}}^{(j)}, \mathbf{y}^{(j)})$ , we have:

$$l_T = l(\mathbf{P}\delta(\mathbf{W}^{(T)})\delta(\mathbf{W}^{(T-1)}) \dots \delta(\mathbf{W}^{(0)})\mathbf{h}^{(0)} + \mathbf{U}^{(0)}\mathbf{x}^{(0)} + \mathbf{b}^{(1)} + \dots + \mathbf{U}^{(T)}\mathbf{x}^{(T)} + \mathbf{b}^{(T)}), \mathbf{y}^{(T)}). \quad (34)$$

Since in RNNs the weight matrices are shared for all time steps:

$$\begin{aligned} \mathbf{W}^{(0)} &= \mathbf{W}^{(1)} = \dots = \mathbf{W}^{(T)} \\ \mathbf{U}^{(0)} &= \mathbf{U}^{(1)} = \dots = \mathbf{U}^{(T)} \\ \mathbf{b}^{(0)} &= \mathbf{b}^{(1)} = \dots = \mathbf{b}^{(T)} \end{aligned}$$

the back propagation needs to be done backward through all time steps ( Back Propagation Through Time-BPTT):

$$\frac{\partial l_T}{\partial \mathbf{W}} = \sum_{t=T}^0 \frac{\partial l_T}{\partial \mathbf{W}^{(t)}}. \quad (35)$$

Each term can be calculated in the same way as MLP:

$$\frac{\partial l_T}{\partial \mathbf{W}^{(t)}} = \frac{\partial l_T}{\partial \mathbf{h}^{(T)}} \left[ \prod_{\tau=T}^{t+1} \text{diag}(\dot{\delta}_\tau) \mathbf{W} \right] \text{diag}(\dot{\delta}_t) \rho(\mathbf{h}^{(t)}, n_h) \quad (36)$$

### 3.2 Gradient Vanishing/Exploding

Note that in Eq 36, the same  $\mathbf{W}$  appears  $T - t$  times in the back propagation. This makes the gradient unstable for RNNs. Consider the linear case,  $\frac{\partial l_T}{\partial \mathbf{W}^{(t)}}$  reduces to:

$$\frac{\partial l_T}{\partial \mathbf{x}^{(t)}} = \frac{\partial l_T}{\partial \mathbf{h}^{(T)}} \left[ \prod_{\tau=T}^t \mathbf{W} \right] \rho(\mathbf{h}^{(t)}, n_h) \quad (37)$$

$$= \frac{\partial l_T}{\partial \mathbf{h}^{(T)}} \mathbf{W}^{T-t} \rho(\mathbf{h}^{(t)}, n_h) \quad (38)$$

- If  $\|\mathbf{W}\| < 1$ ,  $\frac{\partial l_T}{\partial \mathbf{W}^{(t)}}$  vanishes for  $t \ll T$
- If  $\|\mathbf{W}\| > 1$ ,  $\frac{\partial l_T}{\partial \mathbf{W}^{(t)}}$  explodes for  $t \ll T$

Only short term dependencies can be captured!

### Solutions

- Gated architectures: LSTM, GRU
- Constrain  $\|\mathbf{W}\| \approx 1$
- Gradient clipping and etc.

## 4 Convolutional Neural Networks

**Question:** How to deal with image data?

MLP's Problems:

- High dimensional input: RGB image 256x256, input dimension 196608
- Spatial correlations: input features are only locally dependent, no need to capture global correlations.

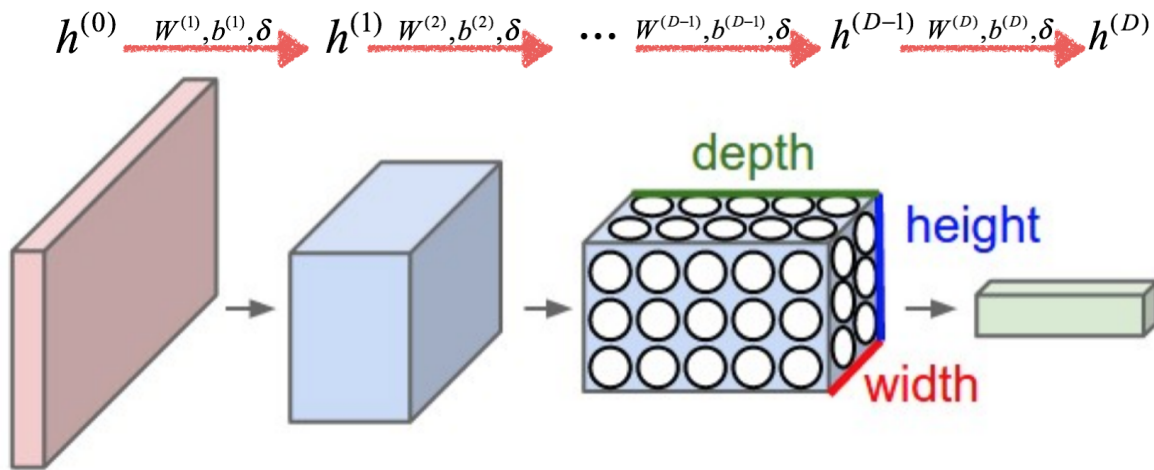


Figure 8: Typical CNN

The most popular model for image data is Convolutional Neural Networks (CNNs).

- Input data and hidden variables are 3rd-order tensors  $shape(x) = (Height, Width, Depth)$ .
- CONV/POOL are two typical kinds of layers

### 4.1 CONV layer

The Conv-layer feature map is defined as:

$$\mathbf{h}_{i,j,k}^{(t)} = \delta \left( \sum_l \left( \mathbf{W}_{::,l,k}^{(t)} * \mathbf{h}_{::,l}^{(t-1)} \right)_{i,j} + \mathbf{b}_k^{(t)} \right)$$

where  $*$  is matrix-matrix conv-operator. For  $\mathbf{P} \in \mathcal{R}^{n_p \times m_p}$  and  $\mathbf{Q} \in \mathcal{R}^{n_q \times m_q}$ , the conv-operator with stride  $s$  is defined as:

$$(\mathbf{P} * \mathbf{Q})_{ij}^{stride=s} = \sum_{a=0}^{n_p} \sum_{b=0}^{m_p} \mathbf{P}_{a,b} \mathbf{Q}_{si-a, sj-b}$$

Note in above equation,  $si - a$  and  $sj - b$  could go out of boundary of  $n_q \times m_q$ . In practice, people usually use three boundary definitions:

- **Hard boundary:** constrain  $0 \leq si - a \leq n_q$  and  $0 \leq sj - b \leq m_q$ .
- **Zero boundary:** filling 0 for out of boundary entrees of  $\mathbf{Q}$
- **Periodic boundary:**  $\mathbf{Q}_{i,j} = \mathbf{Q}_{i \% n_q, j \% m_q}$

Variables of layer  $t$ :

$$\begin{aligned} \mathbf{h}^{(t)} &\in \mathcal{R}^{I^{(t)} \times J^{(t)} \times K^{(t)}} \\ \mathbf{h}^{(t-1)} &\in \mathcal{R}^{I^{(t-1)} \times J^{(t-1)} \times K^{(t-1)}} \\ \mathbf{W}^{(t)} &\in \mathcal{R}^{\tilde{I}^{(t)} \times \tilde{J}^{(t)} \times K^{(t-1)} \times K^{(t)}} \\ \mathbf{b}^{(t)} &\in \mathcal{R}^{K^{(t)}} \end{aligned}$$

where  $I^{(t-1)}, J^{(t-1)}, K^{(t-1)}$  are the height, width and depth of layer  $t - 1$  and the convolutional kernel  $\mathbf{W}^{(t)}$  is used to map

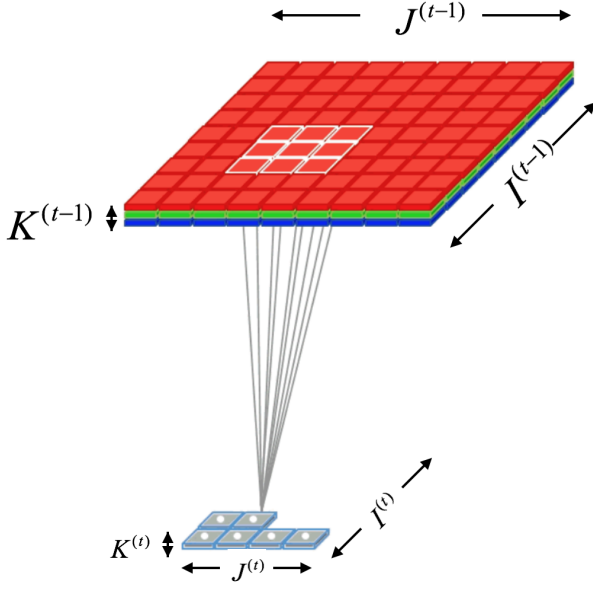


Figure 9: CNN: Conv layer with stride 2.  $\mathbf{h}^{(t-1)} \in \mathcal{R}^{9 \times 9 \times 3}$ ,  $\mathbf{h}^{(t)} \in \mathcal{R}^{4 \times 4 \times 1}$  and  $\mathbf{K} \in \mathcal{R}^{3 \times 3 \times 3 \times 1}$ .

## 4.2 POOL layer

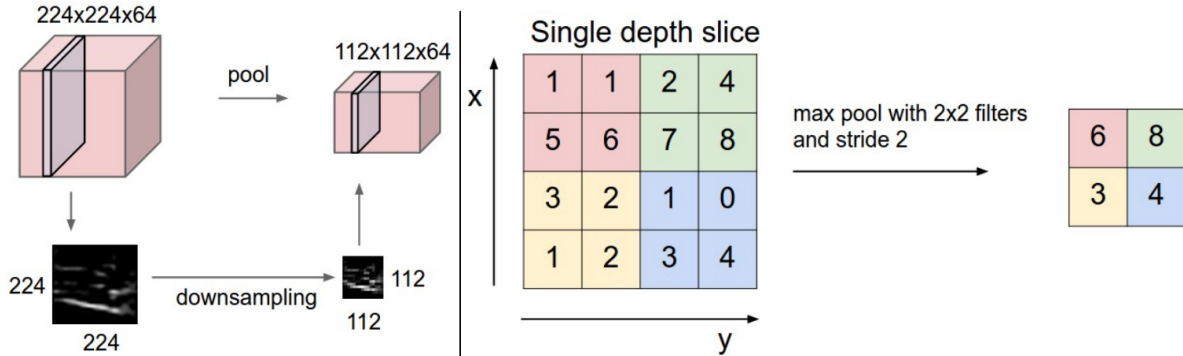


Figure 10: CNN: MAX POOL layer

- To reduce # of parameters, pooling layer is used
- POOL downsamples each depth slice of input feature tensor
- Example: max pooling with stride  $s$ :

$$\mathbf{h}_{i,j,k}^{(t)} = \max(\mathbf{h}_{s_i:(s+1)i, s_j:(s+1)j, k}^{(t-1)})$$

- Max pooling layers have no parameters to learn

## 4.3 Forward Propagation

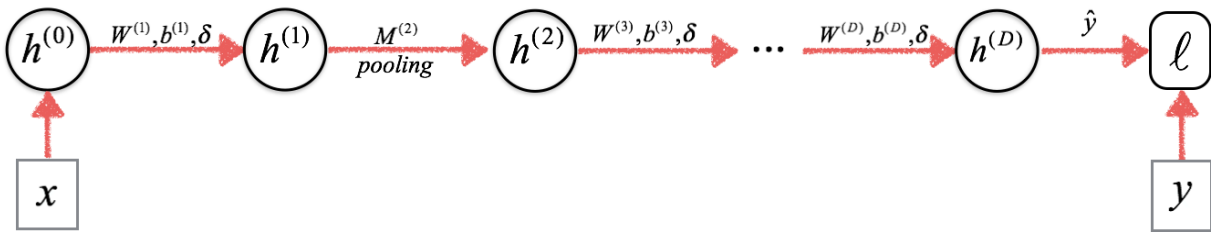


Figure 11: CNN forward propagation

- Evaluate convolutions for each CONV layer

$$\mathbf{h}_{i,j,k}^{(t)} = \delta\left(\sum_l (\mathbf{W}_{::,l,k}^{(t)} * \mathbf{h}_{::,l}^{(t-1)})_{i,j} + \mathbf{b}_k^{(t)}\right)$$

- Evaluate pooling for each POOL layer

$$\mathbf{h}_{i,j,k}^{(t)} = \max(\mathbf{h}_{s_i:(s+1)i-1, s_j:(s+1)j-1, k}^{(t-1)})$$

## 4.4 Back Propagation

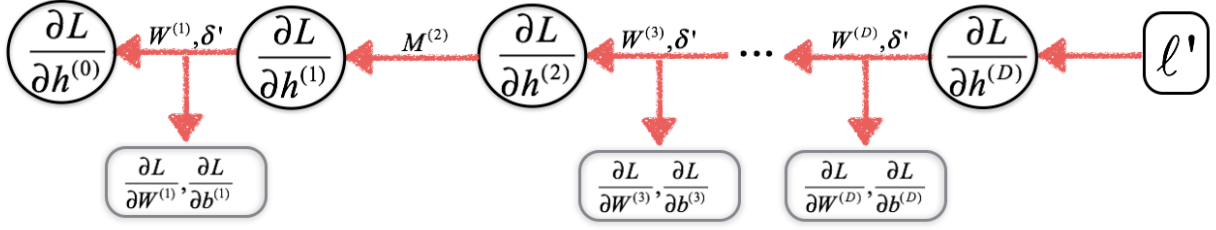


Figure 12: CNN back propagation

- Denoting  $C_{i,j,k} = \delta' \frac{dl}{dh_{i,j,k}^{(t)}}$ , we have:

$$\frac{dl}{dh_{a,b,l}^{(t-1)}} = \sum_{i,j,k} C_{i,j,k} W_{i-a,j-b,l,k}^{(t)} = \sum_k C_{::,k} * W_{::,l,k}^{(t)}$$

$$\frac{dl}{dW_{a,b,l,k}^{(t)}} = \sum_{i,j,k} C_{i,j,k} h_{i-a,j-b,l}^{(t-1)} = \sum_l C_{::,k} * h_{::,l}^{(t-1)}$$

- For POOL layer, the derivative only routing the gradient to the input that had the highest value in the forward pass.