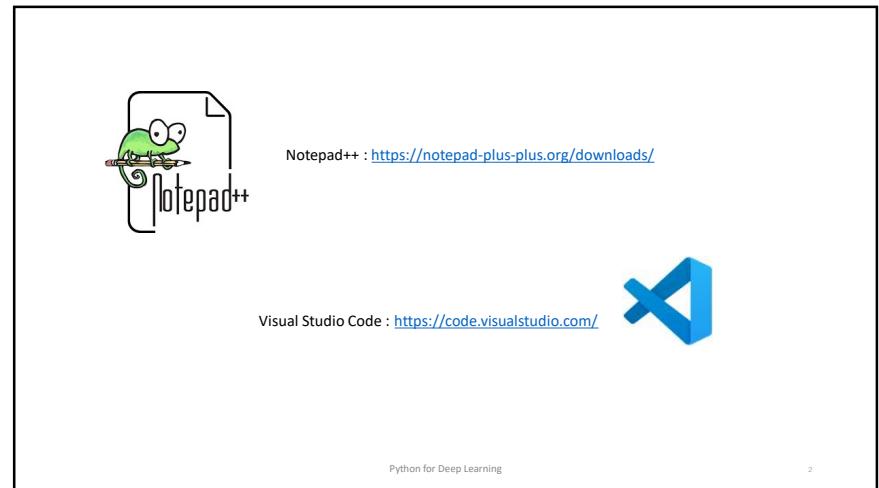


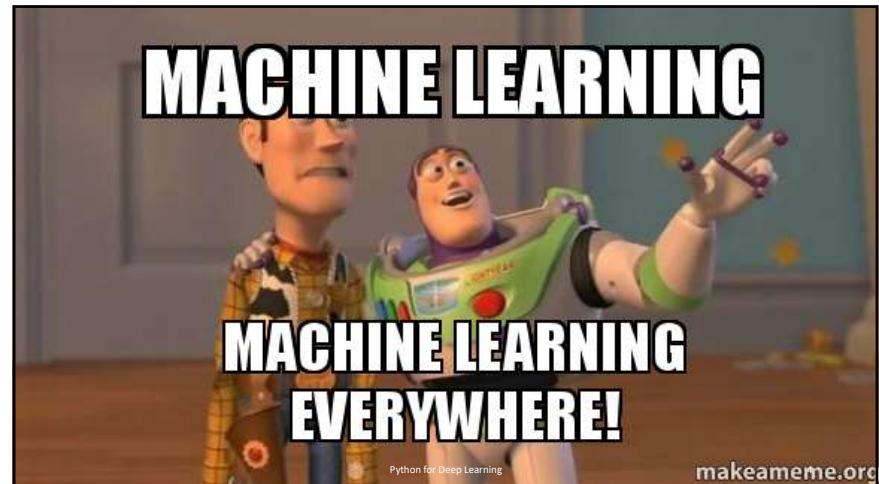
1



2



3



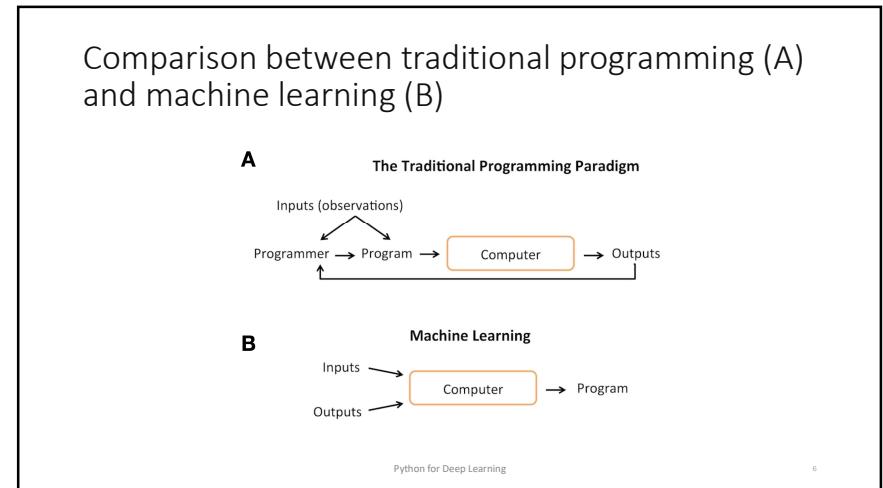
4

What is Machine Learning?

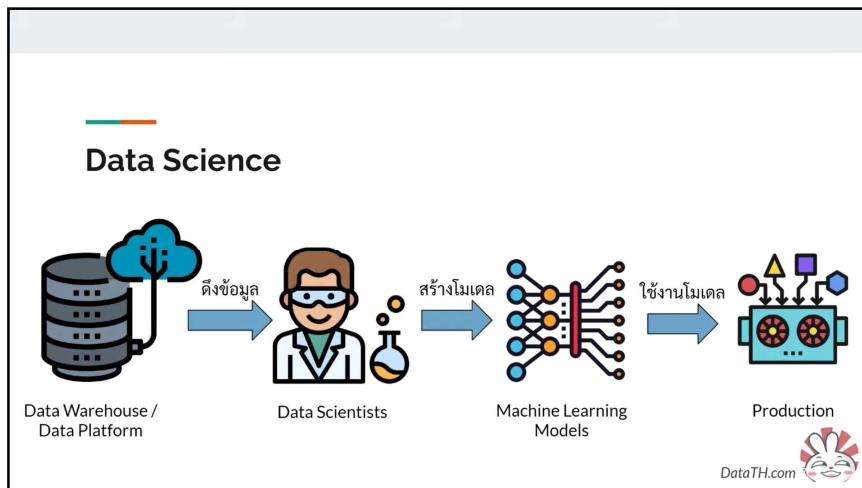
- Machine learning (ML) is the scientific study of algorithms and statistical models that computer systems use to perform a specific task without using explicit instructions, relying on patterns and inference instead.
- https://en.wikipedia.org/wiki/Machine_learning

Python for Deep Learning

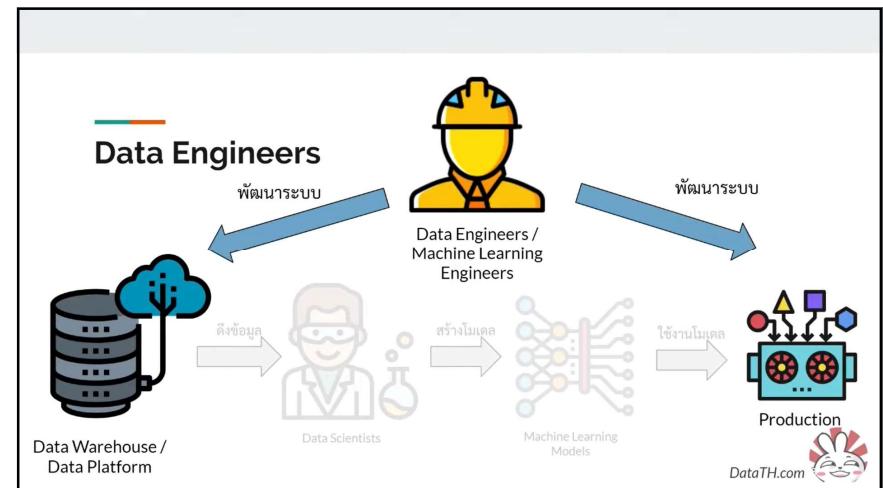
5



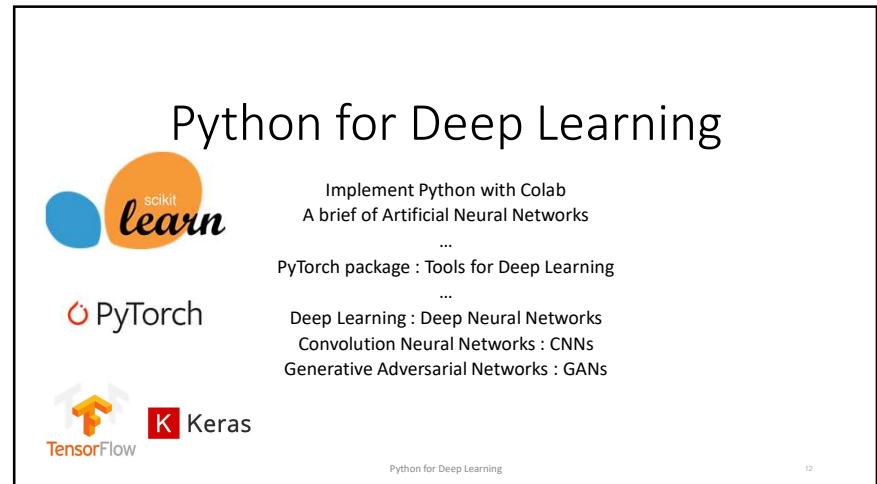
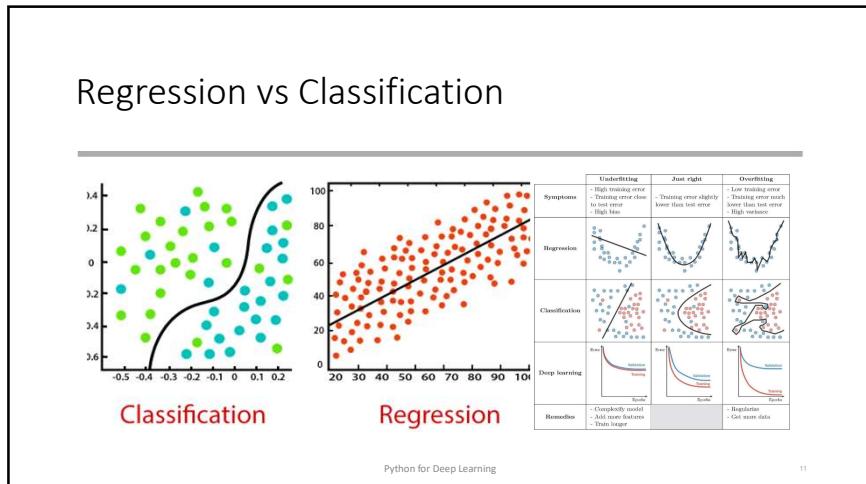
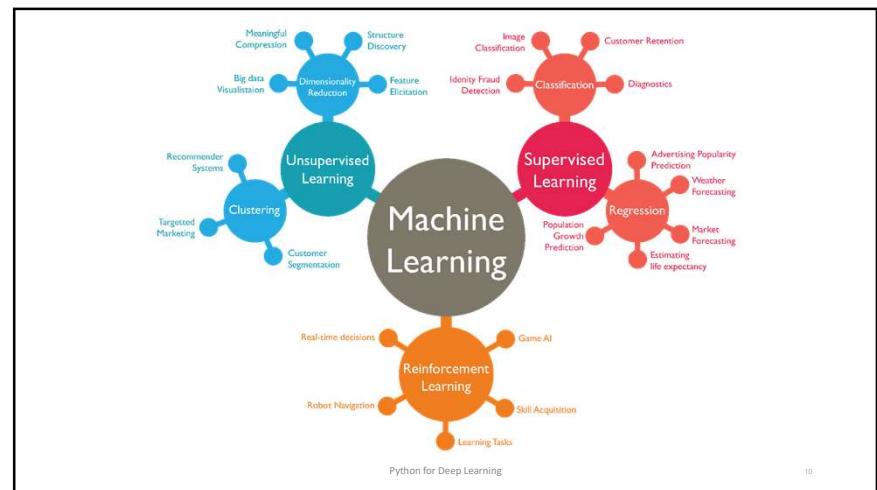
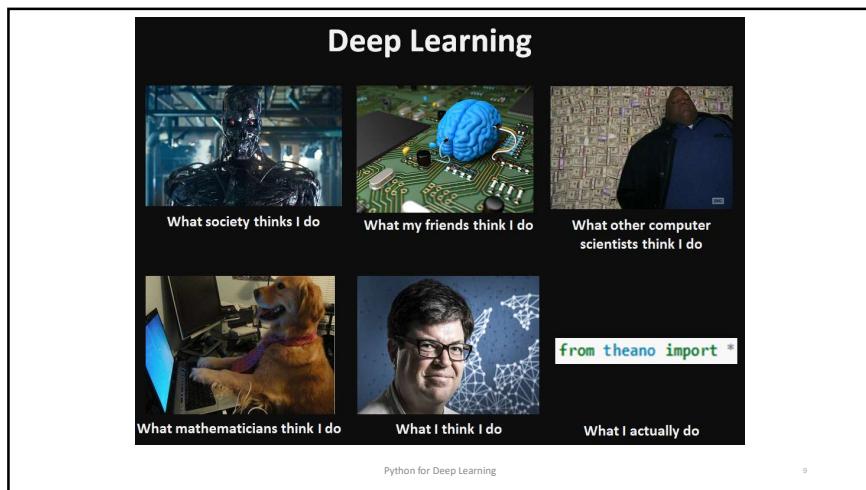
6



7



8



A Brief of Artificial Neural Networks

Perceptron : The First Neural Networks

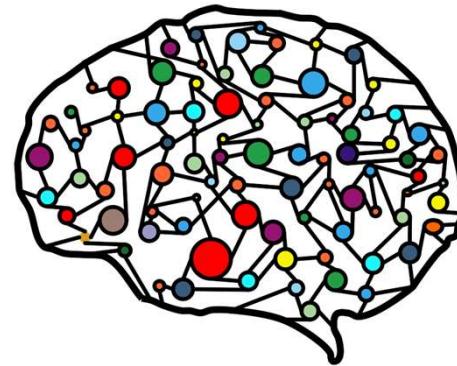
...

Deep Learning : Multi-layer Perceptron / Deep Neural Networks

Python for Deep Learning

13

13



Artificial Neural Networks : ANNs

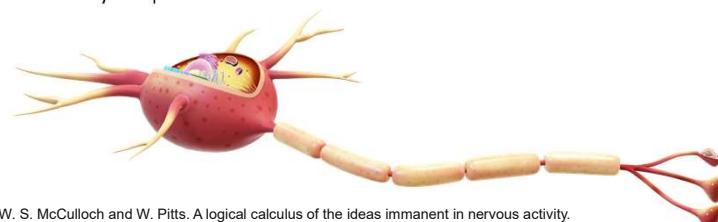
Python for Deep Learning

14

14

A brief of Artificial Neural Networks

- 1943 ANNs began with Warren McCulloch and Walter Pitts [1] who drew an analogy between biological neurons and simple logic gates with binary outputs.



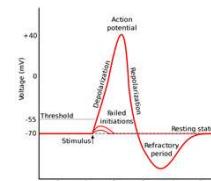
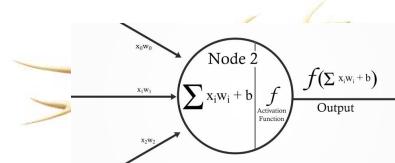
[1] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. The bulletin of mathematical biophysics, 5(4):115–133, 1943.

Python for Deep Learning

15

15

A brief of Artificial Neural Networks



Python for Deep Learning

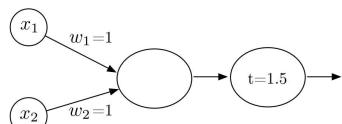
16

16

A brief of Artificial Neural Networks

- 1943 ANNs began with Warren McCulloch and Walter Pitts [1] who drew an analogy between biological neurons and simple logic gates with binary outputs.

x_1	x_2	Out
0	0	0
0	1	0
1	0	0
1	1	1



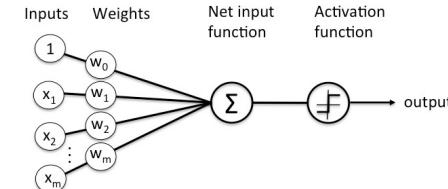
[1] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. The bulletin of mathematical biophysics, 5(4):115–133, 1943.

Python for Deep Learning

17

A brief of Artificial Neural Networks

- 1949 Hebb's rule : “neurons that fire together, wire together” was proposed by Donald O. Hebb.



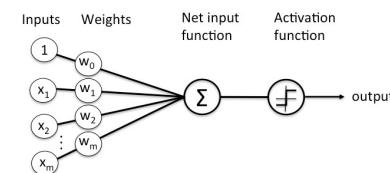
Python for Deep Learning

18

17

A brief of Artificial Neural Networks

- 1957 Frank Rosenblatt [2] published the first concept of the Perceptron learning rule.



[2] F. Rosenblatt. The perceptron, a perceiving and recognizing automaton Project Para. Cornell Aeronautical Laboratory, 1957.

Python for Deep Learning

19

Activation function	Equation	Example	3D Graph
Unit step (Heaviside)	$\phi(z) = \begin{cases} 0, & z < 0, \\ 1, & z \geq 0, \end{cases}$	Perceptron variant	
Sign (Sigmoid)	$\phi(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Linear	$\phi(z) = z$	Adaline, linear regression	
Piece-wise linear	$\phi(z) = \begin{cases} 1, & z \geq 0, \\ 0, & z < 0, \end{cases}$	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, softmax	$\phi(z) = \ln(1 + e^z)$	Multi-layer Neural Networks	

A brief of Artificial Neural Networks

- 1957 Frank Rosenblatt [2] published the first concept of the Perceptron learning rule.

Perceptron Learning Algorithm

$k \leftarrow 1; \mathbf{w}_k \leftarrow \mathbf{0}$.

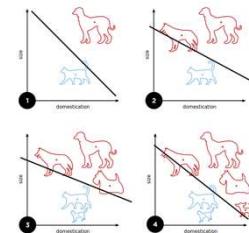
While there exists $i \in \{1, 2, \dots, n\}$ such that $y^i(\mathbf{w}_k \cdot \mathbf{x}^i) \leq 0$:

Pick an arbitrary $j \in \{1, 2, \dots, n\}$ such that $y^j(\mathbf{w}_k \cdot \mathbf{x}^j) \leq 0$.

$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + y^j \mathbf{x}^j$.

$k \leftarrow k + 1$.

Return \mathbf{w}_k .



[2] F. Rosenblatt. The perceptron, a perceiving and recognizing automaton Project Para. Cornell Aeronautical Laboratory, 1957.

Python for Deep Learning

20

19

20

A brief of Artificial Neural Networks

- Rosenblatt's perceptron scheme [2] :

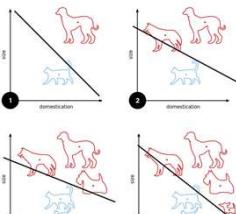
1. Initialize the weights to 0 or small random numbers.

2. For each training sample $x^{(i)}$:

Calculate the output value.

Update the weights.

$$\Delta w_j = \eta(\text{target}^{(i)} - \text{output}^{(i)})x_j^{(i)}$$



[2] F. Rosenblatt. The perceptron, a perceiving and recognizing automaton Project Para. Cornell Aeronautical Laboratory, 1957.

Python for Deep Learning

21

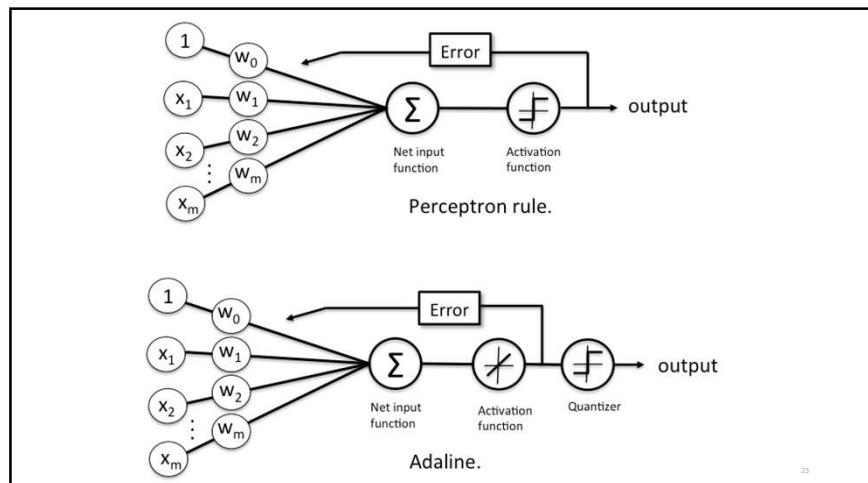
A brief of Artificial Neural Networks

- 1960 Bernard Widrow and Tedd Hoff [3] proposed the idea of the Adaptive Linear Neuron (Adaline).

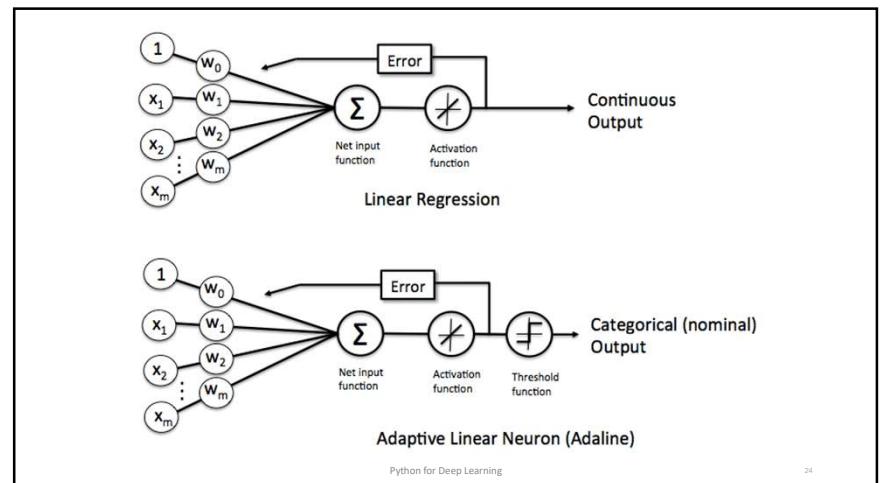
[3] B. Widrow et al. Adaptive "Adaline" neuron using chemical "memistors". Number Technical Report 1553-2. Stanford Electron. Labs., Stanford, CA, October 1960.

Python for Deep Learning

22



23

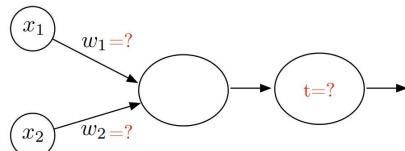


24

A brief of Artificial Neural Networks

- 1967 Marvin L. Minsky and Seymour A. Papert [4] showed that a single-layer perceptron can't implement XOR.

x_1	x_2	Out
0	0	0
0	1	1
1	0	1
1	1	0



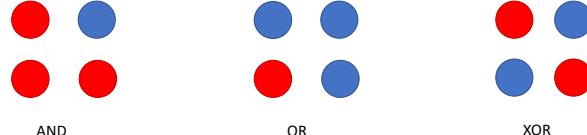
[4] Minski, M. L., & Papert, S. A. (1969). Perceptrons: an introduction to computational geometry. MA: MIT Press, Cambridge.

Python for Deep Learning

25

A brief of Artificial Neural Networks

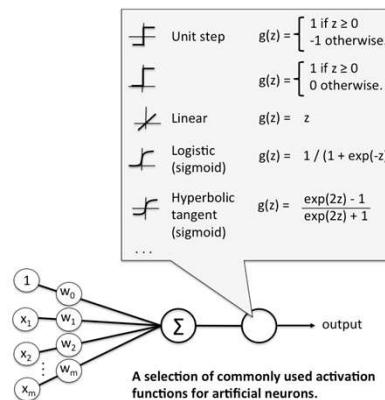
- 1967 Marvin L. Minsky and Seymour A. Papert [4] showed that a single-layer perceptron can't implement XOR.
- Neurons, a dead end? Start of the first "AI Winter"



[4] Minski, M. L., & Papert, S. A. (1969). Perceptrons: an introduction to computational geometry. MA: MIT Press, Cambridge.

Python for Deep Learning

26



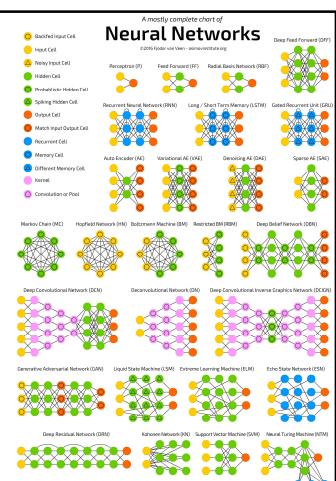
Python for Deep Learning

27

A brief of Artificial Neural Networks

- 1974 Backpropagation was constructed and developed by many researchers, i.e. Paul Werbos, David E. Rumelhart, Geoffrey E. Hinton, Ronald J. Williams
- Who Invented Backpropagation?
- <http://people.idsia.ch/~juergen/who-invented-backpropagation.html>
- 2nd "AI Winter" in the late 1990's and 2000's

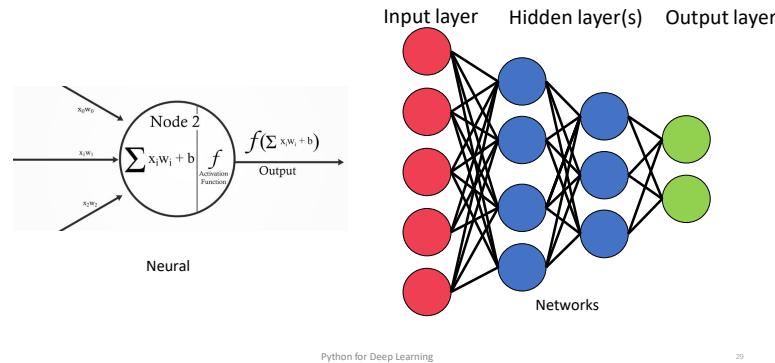
Source : <https://towardsdatascience.com/the-mostly-complete-chart-of-neural-networks-explained-3fb6f2367464>



Python for Deep Learning

28

Artificial Neural Networks



29

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

“On-line” mode

- Initialize $\mathbf{w} := 0^m, b := 0$
- For every training epoch:
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update $\Delta\mathbf{w}, \Delta b$:
 - Update \mathbf{w}, b :

Gradient Descent

- Initialize $\mathbf{w} := 0^m, b := 0$
- For every training epoch:
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - Compute output $\hat{y}^{[i]} = \sigma(\mathbf{w}^T x^{[i]} + b)$
 - Calculate error : $y^{[i]} - \hat{y}^{[i]}$
 - Compute $\nabla_{\mathbf{w}} \mathcal{L} = [\frac{\partial \mathcal{L}}{\partial w_1}, \dots, \frac{\partial \mathcal{L}}{\partial w_m}]^T$ and $\frac{\partial \mathcal{L}}{\partial b}$:
 - For $j \in \{1, 2, \dots, m\}$:
 - $\frac{\partial \mathcal{L}}{\partial w_j} = -\text{error} \times \sigma'(\mathbf{w}^T x^{[i]} + b) x_j^{[i]}$
 - $\frac{\partial \mathcal{L}}{\partial b} = -\text{error} \times \sigma'(\mathbf{w}^T x^{[i]} + b)$
 - Update \mathbf{w}, b :

Python for Deep Learning

30

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

“On-line” mode

- Initialize $\mathbf{w} := 0^m, b := 0$
- For every training epoch:
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update $\Delta\mathbf{w}, \Delta b$:
 - Update \mathbf{w}, b :

Batch mode

- Initialize $\mathbf{w} := 0^m, b := 0$
- For every training epoch:
 - Initialize $\Delta\mathbf{w} := 0, \Delta b := 0$
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update $\Delta\mathbf{w}, \Delta b$:
 - Update \mathbf{w}, b :

Python for Deep Learning

31

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

Batch mode

- Initialize $\mathbf{w} := 0^m, b := 0$
- For every training epoch:
 - Initialize $\Delta\mathbf{w} := 0, \Delta b := 0$
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update $\Delta\mathbf{w}, \Delta b$:
 - Update \mathbf{w}, b :

Batch Gradient Descent

- Initialize $\mathbf{w} := 0^m, b := 0$
- For every training epoch:
 - Initialize $\Delta\mathbf{w} := 0, \Delta b := 0$
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - Compute $\nabla_{\mathbf{w}} \mathcal{L} = [\frac{\partial \mathcal{L}}{\partial w_1}, \dots, \frac{\partial \mathcal{L}}{\partial w_m}]^T$ and $\frac{\partial \mathcal{L}}{\partial b}$:
 - $\Delta\mathbf{w} = \Delta\mathbf{w} + \eta \nabla_{\mathbf{w}} \mathcal{L}, \Delta b = \Delta b + \eta \frac{\partial \mathcal{L}}{\partial b}$
 - Update \mathbf{w}, b :

Python for Deep Learning

32

31

32

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

Gradient Descent ("On-line" mode)

- Initialize $w := 0^m, b := 0$
- For every training epoch:
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - Compute output $\hat{y}^{[i]} = \sigma(w^T x^{[i]} + b)$
 - Calculate error $y^{[i]} - \hat{y}^{[i]}$
 - Compute $\nabla_w \mathcal{L} = \left[\frac{\partial \mathcal{L}}{\partial w_1}, \dots, \frac{\partial \mathcal{L}}{\partial w_m} \right]^T$ and $\frac{\partial \mathcal{L}}{\partial b}$
 - For $j \in \{1, 2, \dots, m\}$:
 - $\frac{\partial \mathcal{L}}{\partial w_j} = -\text{error} \times \sigma'(w^T x^{[i]} + b) x_j^{[i]}$
 - $\frac{\partial \mathcal{L}}{\partial b} = -\text{error} \times \sigma'(w^T x^{[i]} + b)$
 - Update w, b :

$$w := w + \eta \nabla_w \mathcal{L}, b := b + \eta \frac{\partial \mathcal{L}}{\partial b}$$

Python for Deep Learning

Batch Gradient Descent (Batch mode)

- Initialize $w := 0^m, b := 0$
- For every training epoch:
 - Initialize $\Delta w := 0^m, \Delta b := 0$
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - Compute $\nabla_w \mathcal{L} = \left[\frac{\partial \mathcal{L}}{\partial w_1}, \dots, \frac{\partial \mathcal{L}}{\partial w_m} \right]^T$ and $\frac{\partial \mathcal{L}}{\partial b}$
 - Update $\Delta w, \Delta b$:

$$\Delta w := \Delta w + \eta \nabla_w \mathcal{L}, \Delta b := \Delta b + \eta \frac{\partial \mathcal{L}}{\partial b}$$
 - Update w, b :

$$w := w + \Delta w, b := b + \Delta b$$

33

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

"On-line" mode

- Initialize $w := 0^m, b := 0$
- For every training epoch:
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update w, b

Batch mode

- Initialize $w := 0^m, b := 0$
- For every training epoch:
 - Initialize $\Delta w := 0^m, \Delta b := 0$
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update $\Delta w, \Delta b$
 - Update w, b

$$w := w + \Delta w, b := b + \Delta b$$

Python for Deep Learning

34

33

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

"On-line" mode

- Initialize $w := 0^m, b := 0$
- For every training epoch:
 - For every $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update w, b

"On-line" mode II (alternative)

- Initialize $w := 0^m, b := 0$
- For each t iterations:
 - Pick random $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update w, b

Python for Deep Learning

35

35

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

"On-line" mode II (alternative)

- Initialize $w := 0^m, b := 0$
- For each t iterations:
 - Pick random $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update w, b

Stochastic Gradient Descent ("On-line" mode)

- Initialize $w := 0^m, b := 0$
- For each t iterations:
 - Pick random $(x^{[i]}, y^{[i]}) \in D$:
 - (a) Compute output $\hat{y}^{[i]} = \sigma(w^T x^{[i]} + b)$
 - (b) Calculate error $y^{[i]} - \hat{y}^{[i]}$
 - (c) Compute $\nabla_w \mathcal{L} = \left[\frac{\partial \mathcal{L}}{\partial w_1}, \dots, \frac{\partial \mathcal{L}}{\partial w_m} \right]^T$ and $\frac{\partial \mathcal{L}}{\partial b}$
 - Update w, b

$$w := w + \eta \nabla_w \mathcal{L}, b := b + \eta \frac{\partial \mathcal{L}}{\partial b}$$

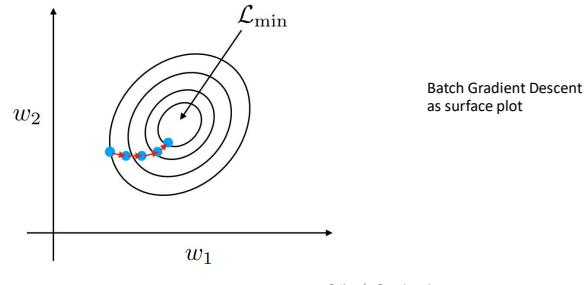
Python for Deep Learning

36

36

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

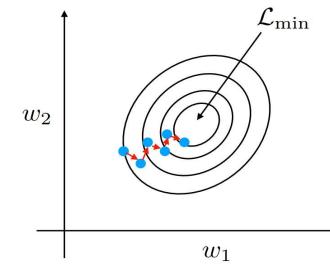


37

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

Stochastic Gradient Descent
as surface plot



Python for Deep Learning

38

General Learning Principle

- Let $D = \{(x^{[i]}, y^{[i]}) \in \mathbb{R}^m \times \{0,1\} : i \in \{1, 2, \dots, n\}\}$

Minibatch mode
(on-line + batch)

Most common used in deep learning!!

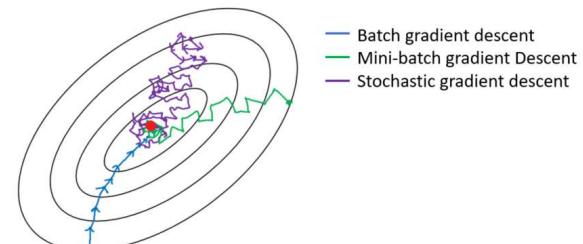
- Initialize $w := 0^m, b := 0$
- For every training epoch:
 - Initialize $\Delta w := 0, \Delta b := 0$
 - For every $\{(x^{[i]}, y^{[i]}), \dots, (x^{[i+k]}, y^{[i+k]})\} \subset D$:
 - (a) Compute output (prediction)
 - (b) Calculate error
 - (c) Update $\Delta w, \Delta b$
 - Update w, b :

$$w := w + \Delta w, b := b + \Delta b$$

Python for Deep Learning

39

General Learning Principle



Source : <https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3>

Python for Deep Learning

40

39

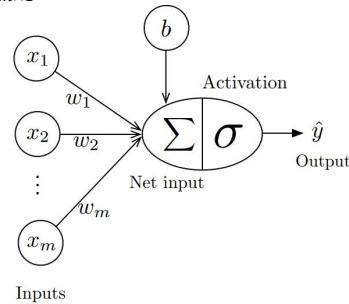
Linear Algebra in DL

- Output $\hat{y} := \sigma(z)$ where $z = x^T w + b, w \in \mathbb{R}^{m \times 1}$

- For 1 example : $x \in \mathbb{R}^{m \times 1}$,
 - $z = x^T w + b$
 - $\hat{y} = \sigma(z)$

- For n examples : $X \in \mathbb{R}^{n \times m}$

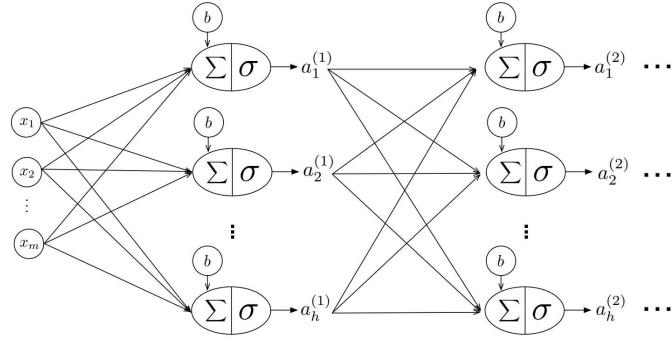
$$\begin{aligned} Xw + b &= \begin{pmatrix} (x^{[1]})^T w + b \\ \vdots \\ (x^{[n]})^T w + b \end{pmatrix} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} = z \in \mathbb{R}^{n \times 1} \\ \hat{y} &= \begin{pmatrix} \sigma(z_1) \\ \vdots \\ \sigma(z_n) \end{pmatrix} = \sigma(z) \end{aligned}$$



Python for Deep Learning

41

Linear Algebra in DL



Python for Deep Learning

42

Linear Algebra in DL (Based on PyTorch)

- Output $a := \sigma(z)$

where $z = xW^T + b$, $W = \begin{bmatrix} w_{1,1} & \cdots & w_{1,m} \\ \vdots & \ddots & \vdots \\ w_{h,1} & \cdots & w_{h,m} \end{bmatrix}$

- For 1 example : $x = [x_1, \dots, x_m] \in \mathbb{R}^{1 \times m}$,

$$\sigma(xW^T + b) = a \in \mathbb{R}^{1 \times h}$$

- For n examples : $X \in \mathbb{R}^{n \times m}$

$$\sigma(XW^T + b) = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} = A \in \mathbb{R}^{n \times h}$$

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

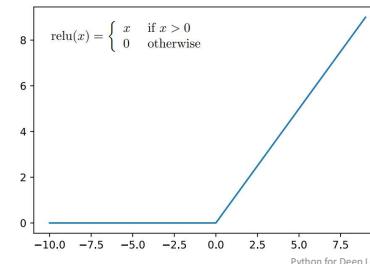
Python for Deep Learning

43

AUTOGRAD: AUTOMATIC DIFFERENTIATION

- Suppose we have the following activation function:

$$a(x, w, b) = \text{relu}(w \cdot x + b)$$



ReLU = Rectified Linear Unit
(prob. the most commonly used activation function in DL)

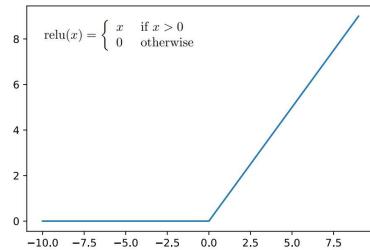
43

44

AUTOGRAD: AUTOMATIC DIFFERENTIATION

- You may note that ($\text{relu} := f$)

$$\cdot f'(x) = \begin{cases} 1, & x > 0; \\ 0, & x < 0; \\ \infty, & x = 0 \end{cases}$$



Python for Deep Learning

45

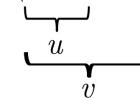
- But in computer science,

$$\cdot f'(x) = \begin{cases} 1, & x > 0; \\ 0, & x \leq 0 \end{cases}$$

45

AUTOGRAD: AUTOMATIC DIFFERENTIATION

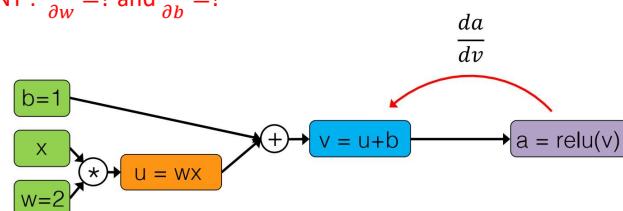
$$a(x, w, b) = \text{relu}(w \cdot x + b)$$

WANT: $\frac{\partial a}{\partial w} = ?$ and $\frac{\partial a}{\partial b} = ?$ 

Python for Deep Learning

46

AUTOGRAD: AUTOMATIC DIFFERENTIATION

WANT: $\frac{\partial a}{\partial w} = ?$ and $\frac{\partial a}{\partial b} = ?$ 

Python for Deep Learning

47

AUTOGRAD: AUTOMATIC DIFFERENTIATION

$$\frac{\partial a}{\partial b} = \frac{\partial v}{\partial b} \frac{\partial a}{\partial v}$$

$$\frac{\partial v}{\partial b}$$

WANT: $\frac{\partial a}{\partial w} = ?$ and $\frac{\partial a}{\partial b} = ?$

$$\frac{\partial a}{\partial w} = \frac{\partial u}{\partial w} \frac{\partial a}{\partial u}$$

$$\frac{\partial u}{\partial w}$$

$$\frac{\partial a}{\partial w} = \frac{\partial u}{\partial w} \frac{\partial v}{\partial u} \frac{\partial a}{\partial v}$$

Python for Deep Learning

48

48

AUTOGRAD: AUTOMATIC DIFFERENTIATION

$$\frac{\partial a}{\partial b} = \frac{\partial v}{\partial b} \frac{\partial a}{\partial v}$$

WANT : $\frac{\partial a}{\partial w} = ?$ and $\frac{\partial a}{\partial b} = ?$

$$\frac{\partial a}{\partial w} = \frac{\partial u}{\partial w} \frac{\partial a}{\partial u}$$

$$= \frac{\partial u}{\partial w} \frac{\partial v}{\partial u} \frac{\partial a}{\partial v}$$

$$= \frac{\partial u}{\partial w} \frac{\partial u}{\partial v} \frac{\partial a}{\partial v}$$

$\text{relu}(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

Python for Deep Learning

49

AUTOGRAD: AUTOMATIC DIFFERENTIATION

$$\frac{\partial a}{\partial b} = \frac{\partial v}{\partial b} \frac{\partial a}{\partial v}$$

WANT : $\frac{\partial a}{\partial w} = ?$ and $\frac{\partial a}{\partial b} = ?$

$$\frac{\partial a}{\partial w} = \frac{\partial u}{\partial w} \frac{\partial a}{\partial u}$$

$$= \frac{\partial u}{\partial w} \frac{\partial v}{\partial u} \frac{\partial a}{\partial v}$$

$$= \frac{\partial u}{\partial w} \frac{\partial v}{\partial w} \frac{\partial a}{\partial v}$$

$\text{relu}(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

Python for Deep Learning

50

49

AUTOGRAD: AUTOMATIC DIFFERENTIATION

$$\frac{\partial a}{\partial b} = \frac{\partial v}{\partial b} \frac{\partial a}{\partial v} = 1$$

WANT : $\frac{\partial a}{\partial w} = ?$ and $\frac{\partial a}{\partial b} = ?$

$$\frac{\partial a}{\partial w} = \frac{\partial u}{\partial w} \frac{\partial a}{\partial u}$$

$$= \frac{\partial u}{\partial w} \frac{\partial v}{\partial u} \frac{\partial a}{\partial v}$$

$$= \frac{\partial u}{\partial w} \frac{\partial v}{\partial w} \frac{\partial a}{\partial v}$$

$$= 3 * 1 * 1 = 3$$

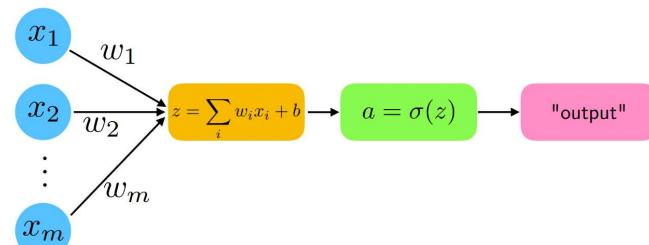
$\text{relu}(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

Python for Deep Learning

51

51

Generic Neuron representation (Recap)

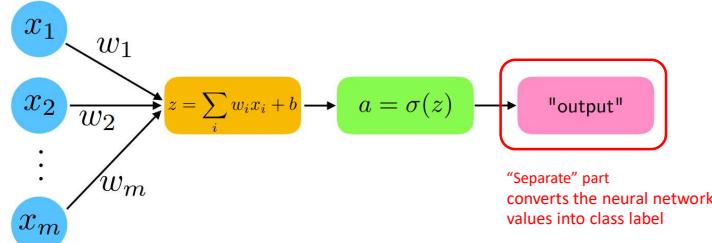


Python for Deep Learning

52

52

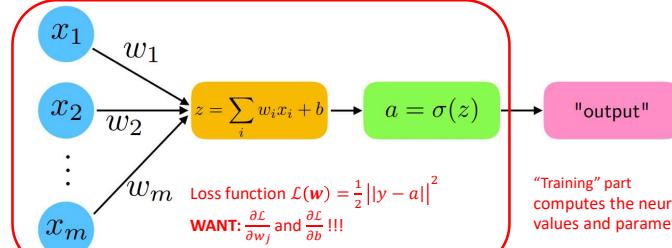
Generic Neuron representation (Recap)



Python for Deep Learning

53

Generic Neuron representation (Recap)

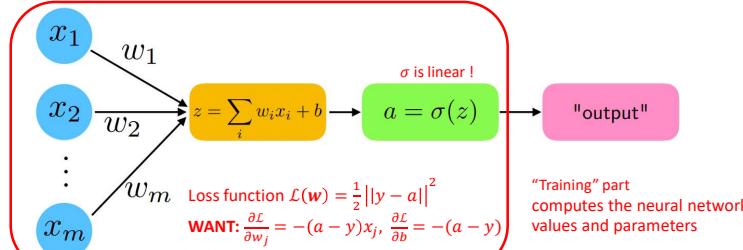


Python for Deep Learning

54

53

Generic Neuron representation (Recap)

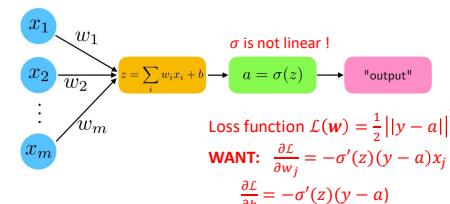


Python for Deep Learning

55

55

Generic Neuron representation (Recap)



Python for Deep Learning

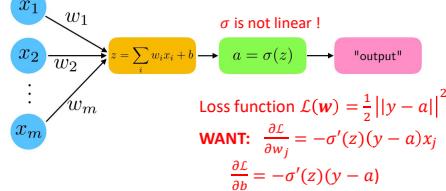
56

56

Activation function	Equation	Example	3D Graph
Unit step (threshold)	$\phi(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}$	Perceptron variant	
Sign (Sigmoid)	$\phi(z) = \begin{cases} -1, & z < 0 \\ 1, & z \geq 0 \end{cases}$	Perceptron variant	
Linear	$\phi(z) = z$	Additive, linear regression	
Piece-wise linear	$\phi(z) = \begin{cases} 1, & z \geq 1 \\ 0, & z < 1 \end{cases}$	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	
Softplus	$\phi(z) = \ln(1 + e^z)$	Multi-layer Neural Networks	

Copyright © Université Paris-Dauphine 2010
http://cermics.enpc.fr/~cermics/teaching/ML/

Logistic Regression



Activation function	Equation	Example	1D Graph
Unit step (Heaviside)	$\phi(z) = \begin{cases} 0, & z < 0, \\ 1, & z \geq 0. \end{cases}$	Perceptron, Ising model	
Sign (Sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Perceptron, Ising model	
Linear	$\phi(z) = z$	Adaline, linear regression	
Piece-wise Linear	$\phi(z) = \begin{cases} 1, & z < \frac{1}{2}, \\ z, & \frac{1}{2} \leq z < 1, \\ 0, & z \geq 1. \end{cases}$	Support vector machine	
Logistic (Sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Logistic, Hyperbolic tangent, Multi-layer Neural Networks	
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	
Receptor, softplus	$\phi(z) = \ln(1 + e^z)$	Multi-layer Neural Networks	

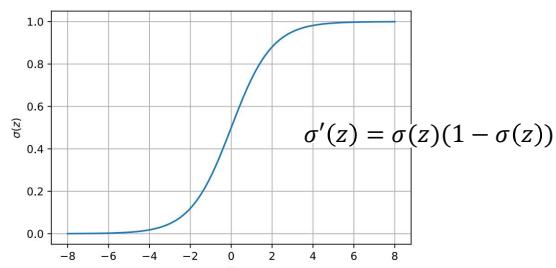
Python for Deep Learning

57

57

Logistic Sigmoid Function

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



Python for Deep Learning

59

59

$$a = \sigma(z) = z \quad \text{Output} = z$$

Linear regression

$$\text{Loss } \mathcal{L}(\mathbf{w}) = \frac{1}{2} \|y - a\|^2$$

$$a = \sigma(z) = z \quad \text{Output} = \text{threshold}(z)$$

Adaline

$$\text{Loss } \mathcal{L}(\mathbf{w}) = \frac{1}{2} \|y - a\|^2$$

$$a = \sigma(z) = \frac{1}{1 + \exp(-z)} \quad \text{Output} = \begin{cases} 1, & \text{if } a > 0.5 \\ 0, & \text{otherwise} \end{cases}$$

Logistic regression

$$\text{Loss } \mathcal{L}(\mathbf{w}) = \frac{1}{2} \|y - a\|^2$$

Python for Deep Learning

58

58

WANT:
 $P(y=0|x) \approx 1$ if $y=0$
 $P(y=1|x) \approx 1$ if $y=1$

- Assume that outputs are independent Bernoulli random variables with success probability of $a = \sigma(\mathbf{w}^T \mathbf{x})$

- The probability mass function P of this distribution over possible outcomes y is

$$P(y|x) = \begin{cases} a, & \text{if } y=1 \\ 1-a, & \text{if } y=0 \end{cases}$$

- Equivalent to,

$$P(y|x) = a^y (1-a)^{(1-y)}$$

Python for Deep Learning

60

60

Logistic Regression

- For n training examples, we want to maximize:

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]})$$

Python for Deep Learning

61

Logistic Regression

- For n training examples, we want to maximize:

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]})$$

- You may know as “Maximum Likelihood Estimation”,

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]}; \mathbf{w}) = \prod_{i=1}^n \sigma(z^{[i]})^{y^{[i]}} (1 - \sigma(z^{[i]}))^{1-y^{[i]}}$$

Python for Deep Learning

62

61

Logistic Regression

- For n training examples, we want to maximize:

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]})$$

- You may know as “Maximum Likelihood Estimation”,

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]}; \mathbf{w}) = \prod_{i=1}^n \sigma(z^{[i]})^{y^{[i]}} (1 - \sigma(z^{[i]}))^{1-y^{[i]}}$$

- In practice, it is easier to take log in to \mathcal{L}_n which called log-likelihood,

$$\ell(\mathbf{w}) = \log(\mathcal{L}(\mathbf{w})) = \sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$$

Python for Deep Learning

63

63

Logistic Regression

In information theory, negative log-likelihood called as binary cross entropy

- For n training examples, we want to maximize:

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]})$$

- You may know as “Maximum Likelihood Estimation”,

$$P(y^{[1]}, y^{[2]}, \dots, y^{[n]} | x^{[1]}, x^{[2]}, \dots, x^{[n]}) = \prod_{i=1}^n P(y^{[i]} | x^{[i]}; \mathbf{w}) = \prod_{i=1}^n \sigma(z^{[i]})^{y^{[i]}} (1 - \sigma(z^{[i]}))^{1-y^{[i]}}$$

- In practice, it is easier to take log in to \mathcal{L}_n which called log-likelihood,

$$\ell(\mathbf{w}) = \log(\mathcal{L}(\mathbf{w})) = \sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$$

- For convenience, we take subtraction to minimize negative log-likelihood instead of maximize log- likelihood,

$$\mathcal{L}(\mathbf{w}) = -\ell(\mathbf{w}) = -\sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$$

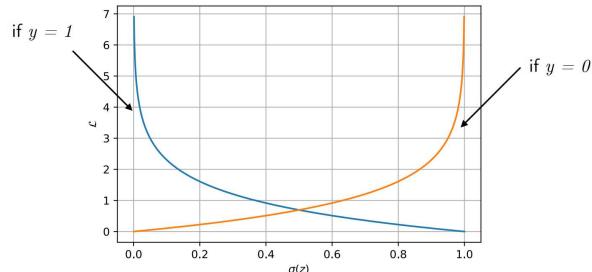
Python for Deep Learning

64

64

Loss for one training example

In information theory, negative log-likelihood called as binary cross entropy



$$\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$$

Python for Deep Learning

65

Logistic Regression : Learning rule

- $\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$
- $\frac{\partial \mathcal{L}}{\partial w_j} = - (y^{[i]} - \sigma(z^{[i]})) x_j$

Python for Deep Learning

66

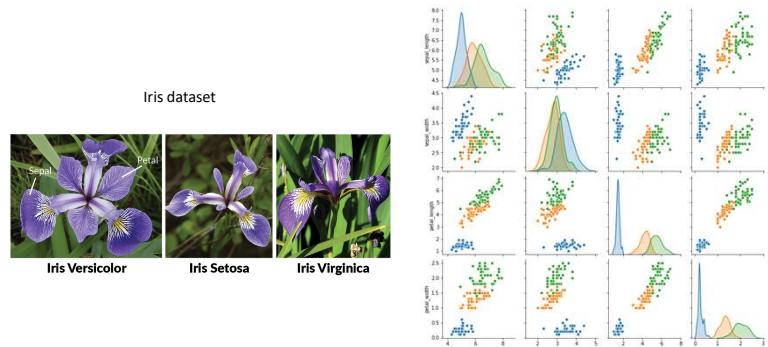
Logistic Regression : Learning rule

- $\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$
- $\frac{\partial \mathcal{L}}{\partial w_j} = - (y^{[i]} - \sigma(z^{[i]})) x_j$
- $w_j = w_j + \eta \frac{\partial \mathcal{L}}{\partial w_j}$
- $w_j = w_j - \eta (y^{[i]} - \sigma(z^{[i]})) x_j \quad \Rightarrow \quad \mathbf{w} = \mathbf{w} + \eta \nabla_{\mathbf{w}} \mathcal{L}$

Python for Deep Learning

67

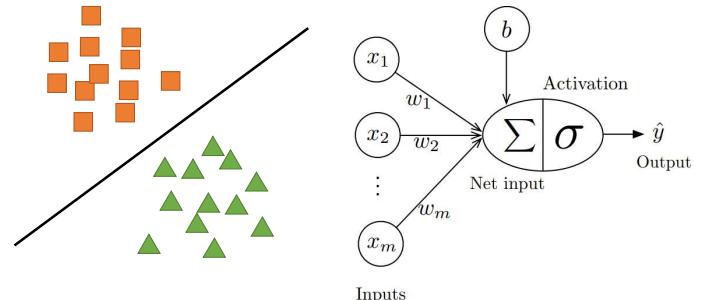
Multi-class Classification



Python for Deep Learning

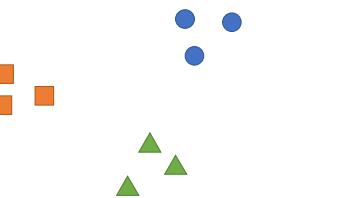
68

Binary-class Classification (Recap)



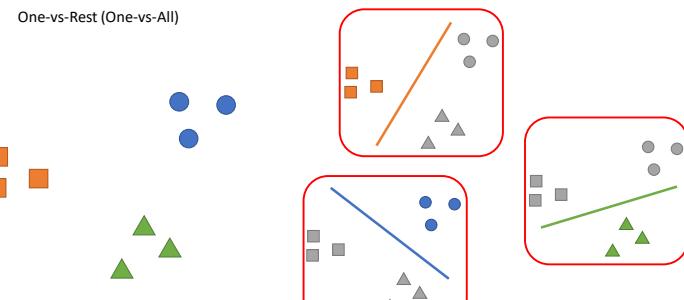
69

Multi-class Classification with Multiple Binary Classifiers



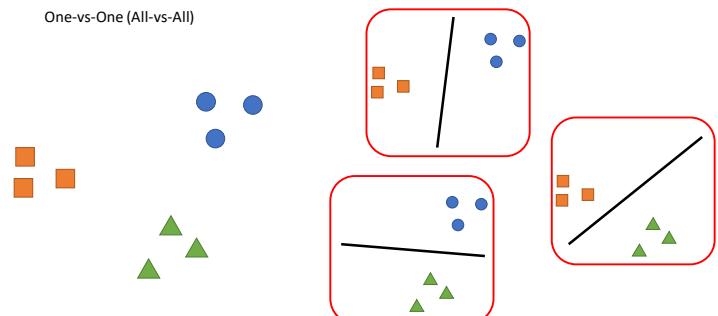
70

Multi-class Classification with Multiple Binary Classifiers



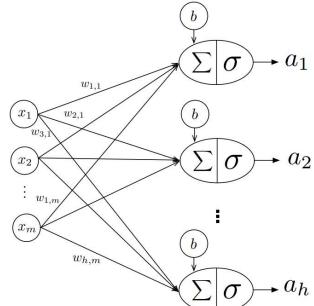
71

Multi-class Classification with Multiple Binary Classifiers



72

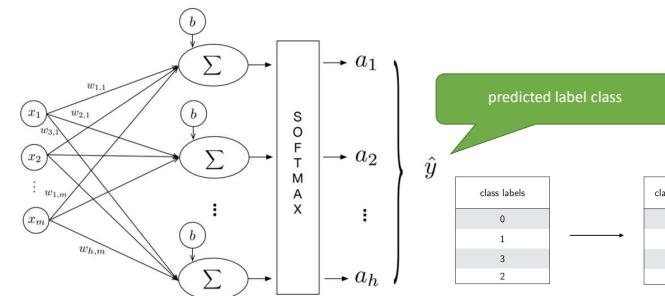
Another approach could be...



Python for Deep Learning

73

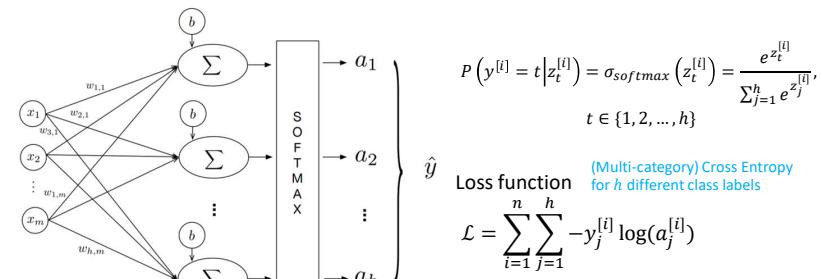
Multinomial Logistic Regression / Softmax Regression



Python for Deep Learning

74

Multinomial Logistic Regression / Softmax Regression



Python for Deep Learning

75

Logistic Regression vs Softmax Regression

	Logistic Regression	Softmax Regression
Activation function : σ	Sigmoid : $\sigma_{sigmoid} = \frac{e^z}{1 + e^z}$	Softmax : $\sigma_{softmax} = \frac{e^{z_t}}{\sum_{j=1}^h e^{z_j}}$
Loss function : \mathcal{L}	Binary cross entropy : $\mathcal{L} = -\sum_{i=1}^n y^{[i]} \log(\sigma(z^{[i]})) + (1 - y^{[i]}) \log(1 - \sigma(z^{[i]}))$	(Multi-category) Cross Entropy : $\mathcal{L} = \sum_{i=1}^n \sum_{j=1}^h -y_j^{[i]} \log(a_j^{[i]})$

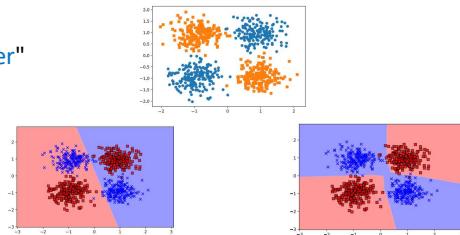
Python for Deep Learning

76

Multi-layer Perceptron : MLP

- 1967 Marvin L. Minsky and Seymour A. Papert [4] showed that a single-layer perceptron can't implement XOR.
- Neurons, a dead end?

Start of the first "AI Winter"

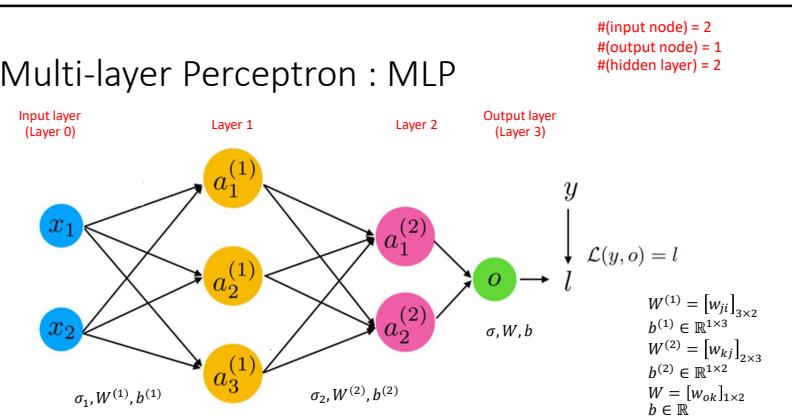


[4] Minsky, M. L., & Papert, S. A. (1969). Perceptrons: an introduction to computational geometry. MA: MIT Press, Cambridge.

Python for Deep Learning

77

Multi-layer Perceptron : MLP

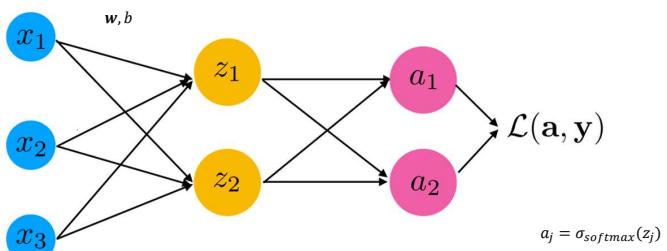


Python for Deep Learning

78

Softmax Regression structure (recap)

#(input node) = 3
#(output node) = 2
#(hidden layer) = 0



Python for Deep Learning

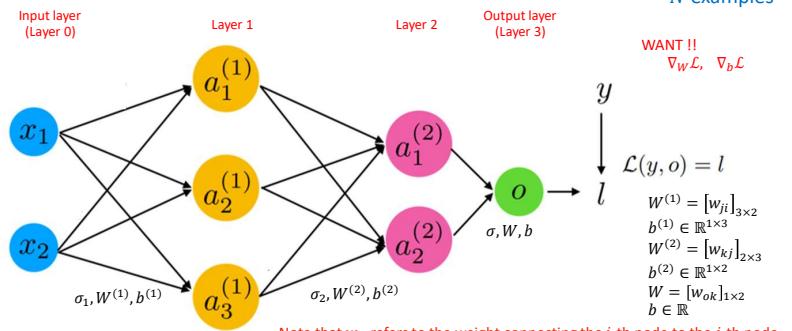
79

Multi-layer Perceptron : MLP

$$\mathcal{L}(\mathbf{y}, \mathbf{o}) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - o^{(n)})^2$$

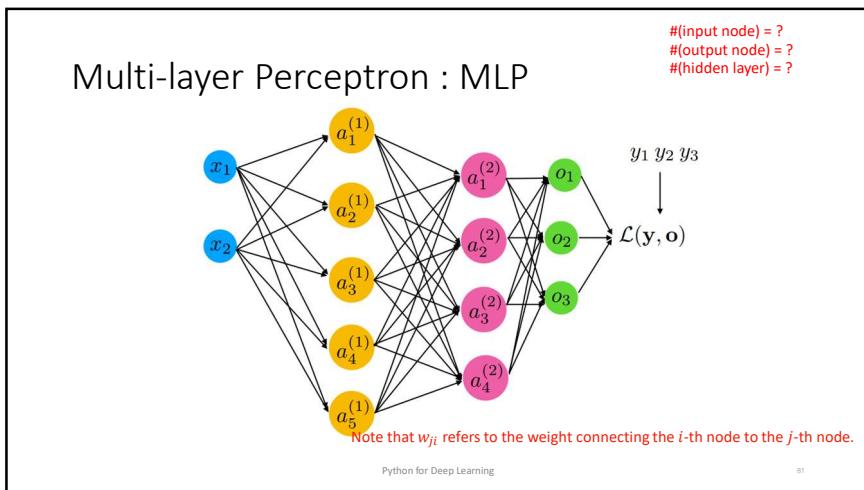
N examples

WANT !!
 $\nabla_w \mathcal{L}, \nabla_b \mathcal{L}$

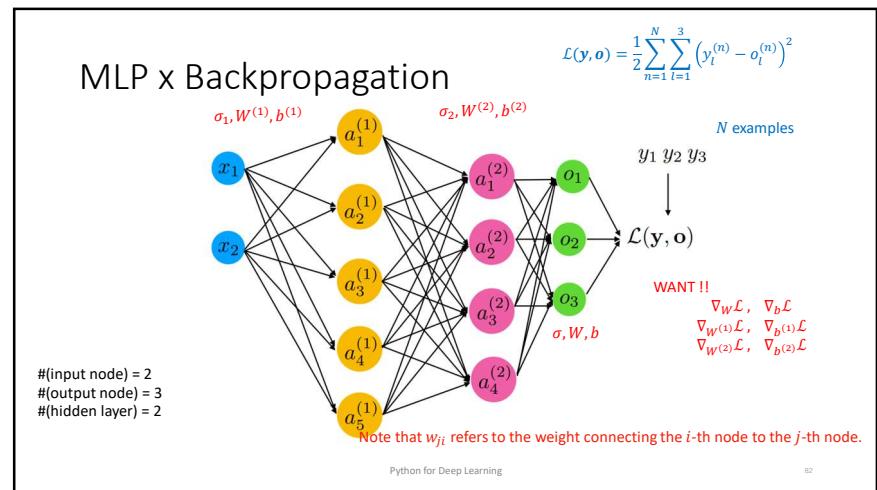


Python for Deep Learning

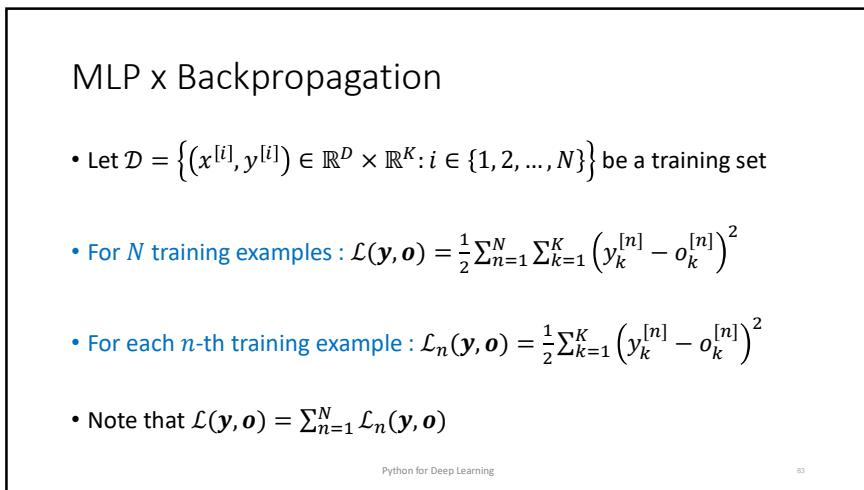
80



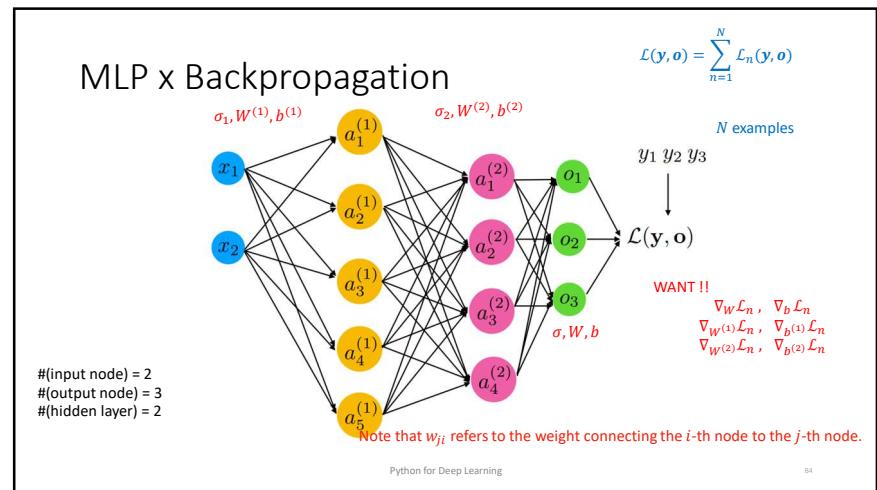
81



82



83



84

MLP x Backpropagation (Simple example)

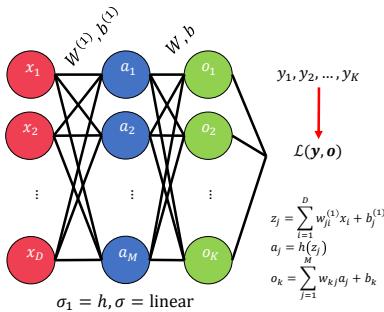
```
#(training example) = 1
#(input node) = D
#(output node) = K
#(hidden layer) = 1
#(hidden node) = M
```

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

85

MLP x Backpropagation (Simple example)



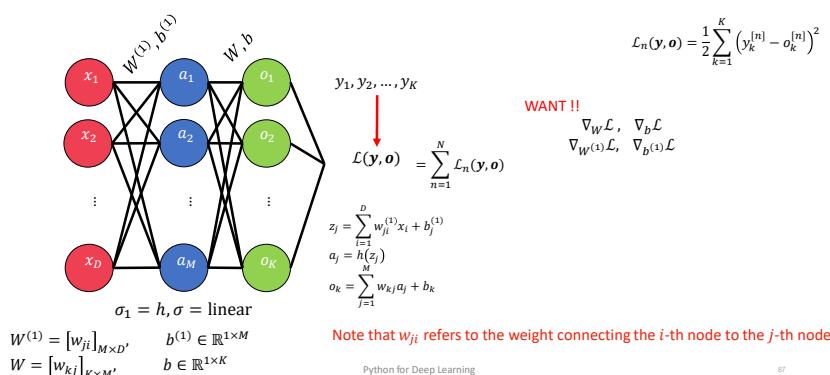
```
#(training example) = 1
#(input node) = D
#(output node) = K
#(hidden layer) = 1
#(hidden node) = M
```

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

86

MLP x Backpropagation (Simple example)



$$\mathcal{L}_n(\mathbf{y}, \mathbf{o}) = \frac{1}{2} \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]})^2$$

WANT !!

$$\nabla_W \mathcal{L}, \nabla_b \mathcal{L}$$

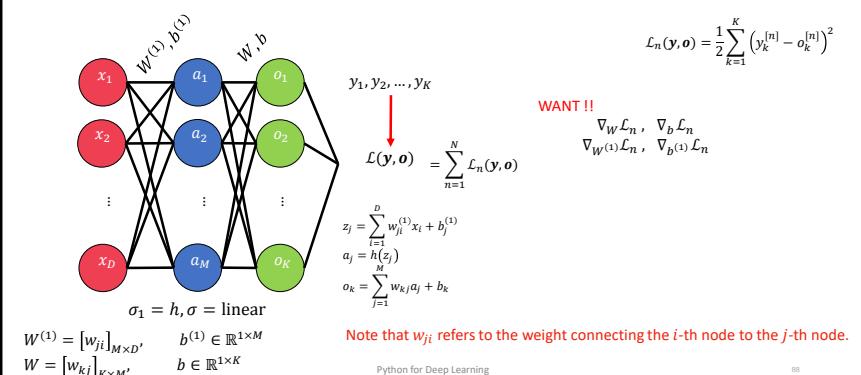
$$\nabla_{W^{(1)}} \mathcal{L}, \nabla_{b^{(1)}} \mathcal{L}$$

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

87

MLP x Backpropagation (Simple example)



$$\mathcal{L}_n(\mathbf{y}, \mathbf{o}) = \frac{1}{2} \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]})^2$$

WANT !!

$$\nabla_W \mathcal{L}_n, \nabla_b \mathcal{L}_n$$

$$\nabla_{W^{(1)}} \mathcal{L}_n, \nabla_{b^{(1)}} \mathcal{L}_n$$

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

88

MLP x Backpropagation (Simple example)

$\sigma_1 = h, \sigma = \text{linear}$

$W^{(1)} = [w_{ji}]_{M \times D}, \quad b^{(1)} \in \mathbb{R}^{1 \times M}$

$W = [w_{kj}]_{K \times M}, \quad b \in \mathbb{R}^{1 \times K}$

$L_n(\mathbf{y}, \mathbf{o}) = \frac{1}{2} \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]})^2$

WANT !!

$$\nabla_W L_n = \left[\frac{\partial L_n}{\partial w_{1,1}}, \dots, \frac{\partial L_n}{\partial w_{K,1}} \right]^T,$$

$$\nabla_b L_n = \left[\frac{\partial L_n}{\partial b_1}, \dots, \frac{\partial L_n}{\partial b_K} \right]^T,$$

$$\nabla_{W^{(1)}} L_n = \left[\frac{\partial L_n}{\partial w_{1,1}^{(1)}}, \dots, \frac{\partial L_n}{\partial w_{M,1}^{(1)}} \right]^T,$$

$$\nabla_{b^{(1)}} L_n = \left[\frac{\partial L_n}{\partial b_1^{(1)}}, \dots, \frac{\partial L_n}{\partial b_M^{(1)}} \right]^T$$

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

89

MLP x Backpropagation (Simple example)

$\sigma_1 = h, \sigma = \text{linear}$

$W^{(1)} = [w_{ji}]_{M \times D}, \quad b^{(1)} \in \mathbb{R}^{1 \times M}$

$W = [w_{kj}]_{K \times M}, \quad b \in \mathbb{R}^{1 \times K}$

$L_n(\mathbf{y}, \mathbf{o}) = \frac{1}{2} \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]})^2$

WANT !!

$$\frac{\partial L_n}{\partial w_{kj}} = ???$$

$$\frac{\partial L_n}{\partial b_k} = ???$$

$$\frac{\partial L_n}{\partial w_{ji}^{(1)}} = ???$$

$$\frac{\partial L_n}{\partial b_m^{(1)}} = ???$$

I KNOW

THE CHAIN RULE

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

90

MLP x Backpropagation (Simple example)

$\sigma_1 = h, \sigma = \text{linear}$

$W^{(1)} = [w_{ji}]_{M \times D}, \quad b^{(1)} \in \mathbb{R}^{1 \times M}$

$W = [w_{kj}]_{K \times M}, \quad b \in \mathbb{R}^{1 \times K}$

$L_n(\mathbf{y}, \mathbf{o}) = \frac{1}{2} \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]})^2$

WANT !!

$$\frac{\partial L_n}{\partial w_{kj}} = -(y_k^{[n]} - o_k^{[n]}) a_j$$

$$\frac{\partial L_n}{\partial b_k} = -(y_k^{[n]} - o_k^{[n]})$$

$$\frac{\partial L_n}{\partial w_{ji}^{(1)}} = - \left(h'(z_j) \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]}) w_{kj} \right) x_i$$

$$\frac{\partial L_n}{\partial b_j^{(1)}} = - \left(h'(z_j) \sum_{k=1}^K (y_k^{[n]} - o_k^{[n]}) w_{kj} \right)$$

Note that w_{ji} refers to the weight connecting the i -th node to the j -th node.

Python for Deep Learning

91

White vs Deep Architectures ?

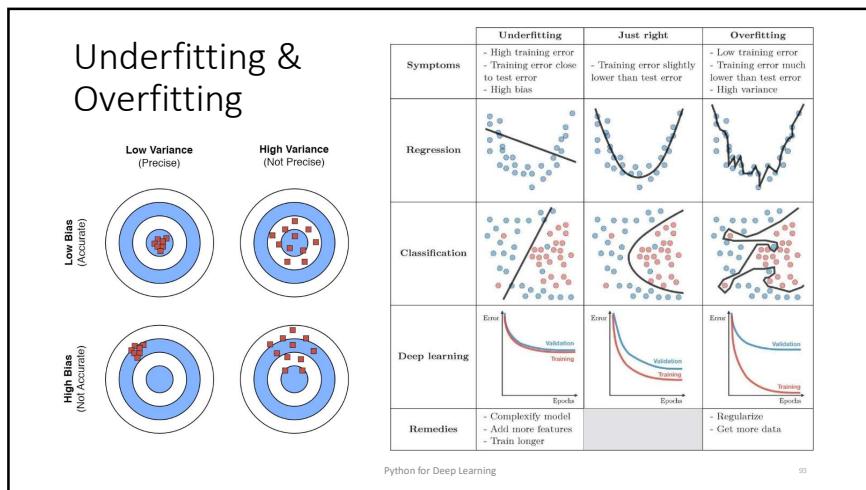
Breadth vs Depth

Parameters = $(4 \times 5) + (3 \times 4) + (2 \times 3) = 41$

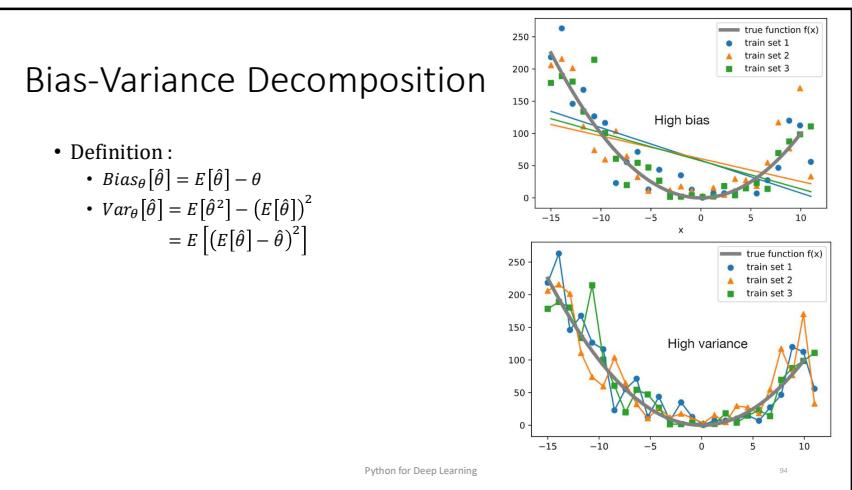
Parameters = $(7 \times 5) + (2 \times 7) + 2 = 51$

Python for Deep Learning

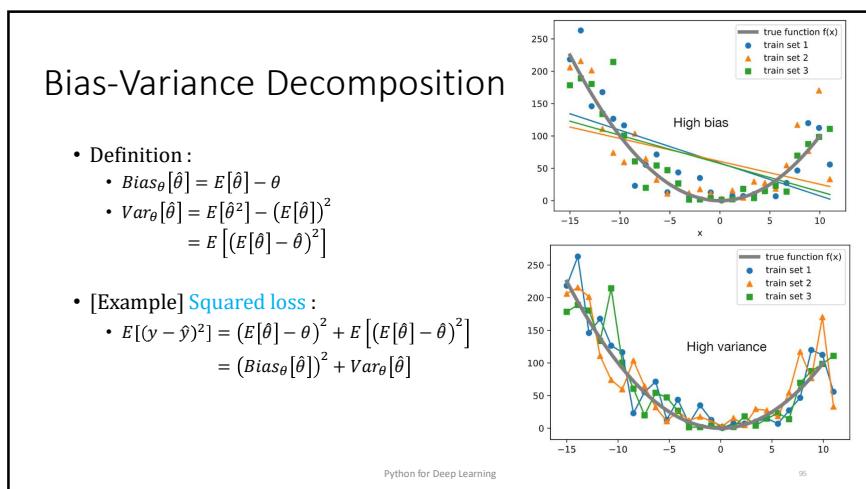
92



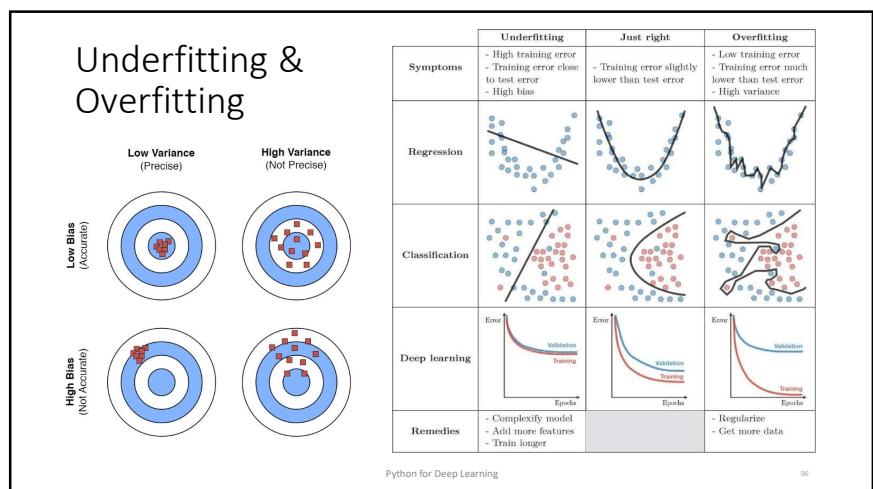
93



94



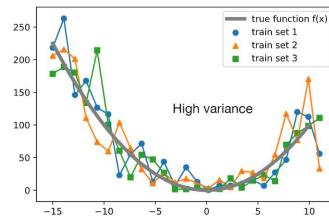
95



96

How to reduce overfitting

- Early stopping
- Regularization
- Dropout

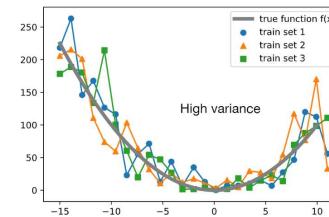


Python for Deep Learning

97

How to reduce overfitting

- Early stopping
- Regularization
- Dropout essentially decreases the size of the no. of parameters to be accounted for during the process of training/learning



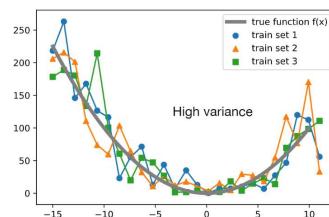
Python for Deep Learning

98

97

How to reduce overfitting

- Early stopping
- **Regularization**
- Dropout



Python for Deep Learning

99

99

How to reduce overfitting

- Early stopping
- Regularization
- Dropout

- L_1 -regularization \Rightarrow LASSO
$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(y^{[i]}, o^{[i]}) + \lambda \|\mathbf{w}\|_1$$
- L_2 -regularization \Rightarrow Ridge regression (Tikhonov regularization)
$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(y^{[i]}, o^{[i]}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$
- $\alpha L_1 + (1 - \alpha) L_2$ -regularization \Rightarrow Elastic-net
$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(y^{[i]}, o^{[i]}) + \lambda \left[\alpha \|\mathbf{w}\|_1 + (1 - \alpha) \frac{1}{2} \|\mathbf{w}\|_2^2 \right]$$

Python for Deep Learning

100

100

Computer can see images as ...

2
4
6
8

5 10 15

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 0 & 4 & 4 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \\ 0 & 8 & 8 & 0 & 0 & 0 & 4 & 4 & 0 \end{bmatrix}$$

Python for Deep Learning 101

101

Computer can see images as ...

2
4
6
8

5 10 15

$$x(:,:,1) = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$x(:,:,2) = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$x(:,:,3) = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Python for Deep Learning 102

102

Why image classification is Hard

Image Source: twitter.com/%2FCats&psig=AOvVaw30_o-PCM-K21DiMAQimQ4&ust=155388775741551

Image Source: https://www.123f.com/photo_76714328_side-view-of-tabby-cat-face-over-white.html

Python for Deep Learning 103

103

Convolution Neural Networks : CNNs

INPUT CONVOLUTION + RELU POOLING CONVOLUTION + RELU POOLING FEATURE LEARNING FLATTEN FULLY CONNECTED CLASSIFICATION

CAR TRUCK VAN
BICYCLE

Python for Deep Learning 104

104

IMAGENET

Explore Download Challenges Publications Updates About
Not logged in. Login | Signup

14,197,122 images, 21841 synsets indexed

ImageNet is an image database organized according to the WordNet hierarchy (currently only the nouns), in which each node of the hierarchy is depicted by hundreds and thousands of images. Currently we have an average of over five hundred images per node. We hope ImageNet will become a useful resource for researchers, educators, students and all of you who share our passion for pictures.
Click here to learn more about ImageNet. Click here to join the ImageNet mailing list.

What do these images have in common? Find out!

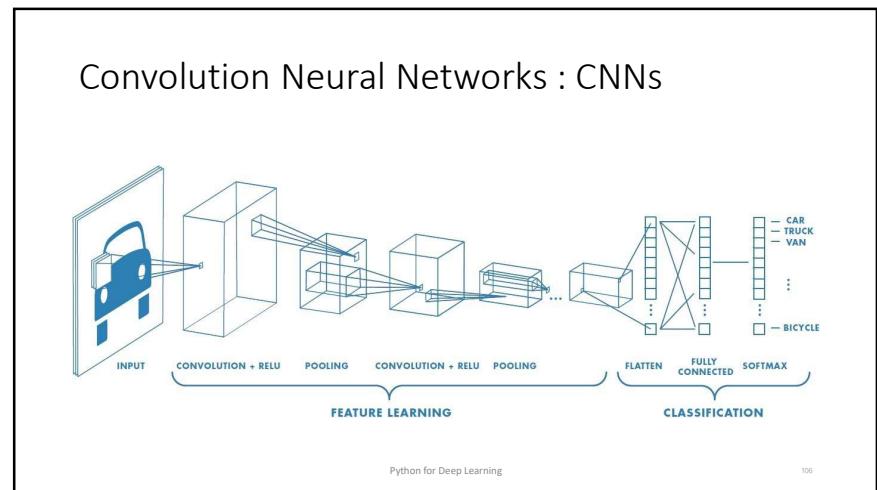
Research updates on improving ImageNet data

© 2016 Stanford Vision Lab, Stanford University, Princeton University support@image-net.org Copyright infringement

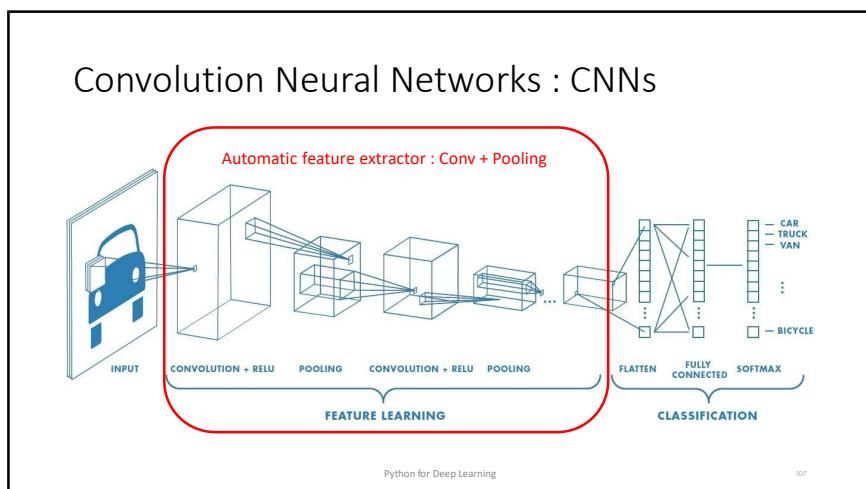
ImageNet Large Scale Visual Recognition Challenge : ILSVRC

Python for Deep Learning 105

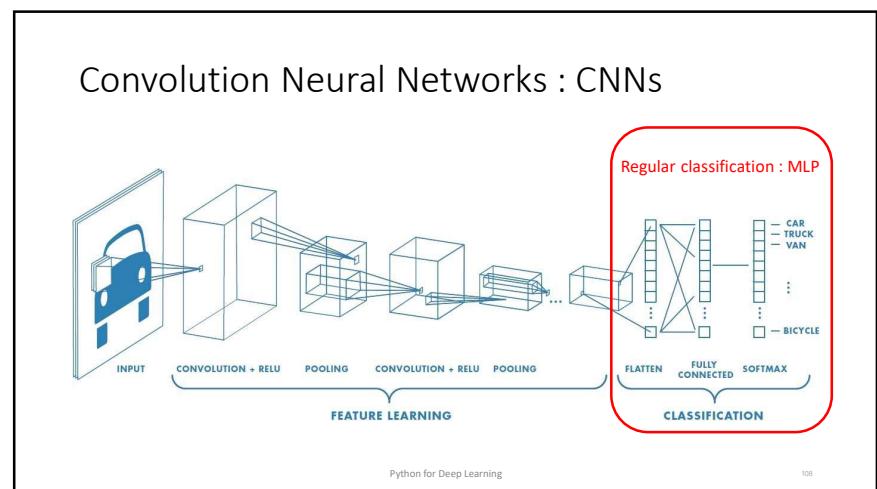
105



106



107



108

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks

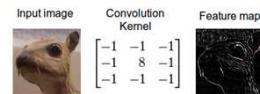
Python for Deep Learning

109

109

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks



$$y_j = \sum_{i=0}^8 w_i x_k$$

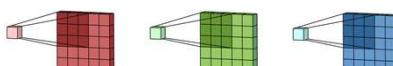
Python for Deep Learning

110

110

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks



Python for Deep Learning

111

111

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks



Python for Deep Learning

112

112

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks

The diagram shows four stages of a CNN:

- Convolution:** An input layer (3x3 grid) is processed by a convolution step to produce an output layer (2x2 grid).
- Stride/Padding:** Two examples are shown:
 - No padding, stride=1:** Input (3x3) becomes output (2x2).
 - padding=2, stride=1:** Input (3x3) becomes output (2x2), with a larger intermediate feature map (5x5) shown.
- Pooling:** An input layer (3x3 grid) is processed by a pooling step to produce an output layer (2x2 grid).
- Full connected networks:** An input layer (2x2 grid) is fully connected to a single output node.

Python for Deep Learning 113

113

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks

The diagram shows four stages of a CNN:

- Convolution:** An input layer (3x3 grid) is processed by a convolution step to produce an output layer (2x2 grid).
- Stride/Padding:** Two examples are shown:
 - No padding, stride=1:** Input (3x3) becomes output (2x2).
 - padding=2, stride=1:** Input (3x3) becomes output (2x2), with a larger intermediate feature map (5x5) shown.
- Pooling:** An input layer (3x3 grid) is processed by a pooling step to produce an output layer (2x2 grid).
- Full connected networks:** An input layer (2x2 grid) is fully connected to a single output node.

Python for Deep Learning 114

114

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks

The diagram shows four stages of a CNN:

- Convolution:** An input layer (3x3 grid) is processed by a convolution step to produce an output layer (2x2 grid).
- Stride/Padding:** Two examples are shown:
 - No padding, stride=1:** Input (3x3) becomes output (2x2). Formula: $(4 - 3 + 2*0)/1 + 1 = 2$
 - padding=2, stride=1:** Input (3x3) becomes output (2x2). Formula: $(5 - 4 + 2*2)/1 + 1 = 6$
- Pooling:** An input layer (3x3 grid) is processed by a pooling step to produce an output layer (2x2 grid).
 - No padding, stride=2:** Input (3x3) becomes output (2x2). Formula: $(5 - 3 + 2*0)/2 + 1 = 2$
- Full connected networks:** An input layer (2x2 grid) is fully connected to a single output node.

Python for Deep Learning 115

115

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks

Feature Map

6	6	6	6
4	5	5	4
2	4	4	2
2	4	4	2

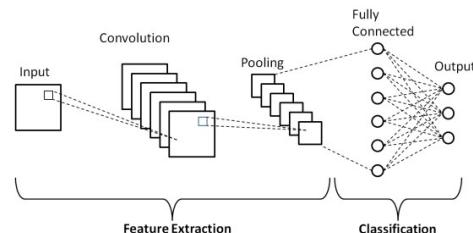
Max Pooling Average Pooling Sum Pooling

Python for Deep Learning 116

116

CNNs in 4 steps

- Convolution
- Stride/Padding
- Pooling
- Full connected networks



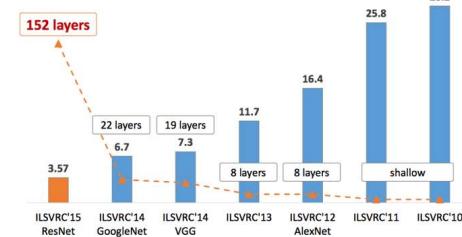
Python for Deep Learning

117

117

Various Types of CNNs

- LeNet
- AlexNet
- VGG
- GoogleNet/Inception
- ResNets

Source : <https://towardsdatascience.com/varioustypes-of-convolutional-neural-network-8b00c9a08a1b>

Python for Deep Learning

118

118

Various Types of CNNs

- LeNet
- AlexNet
- VGG
- GoogleNet
- ResNets

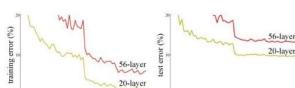


Figure 1: Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer "plain" networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

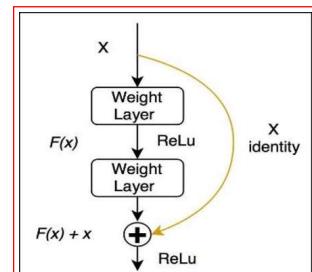
Source : <https://towardsdatascience.com/varioustypes-of-convolutional-neural-network-8b00c9a08a1b>

Python for Deep Learning

119

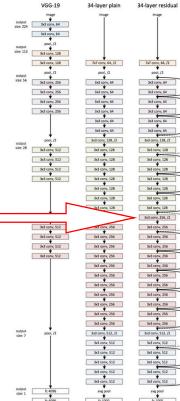
119

ResNets : Residual Networks

Paper : <https://arxiv.org/abs/1512.03385>

Python for Deep Learning

120



120

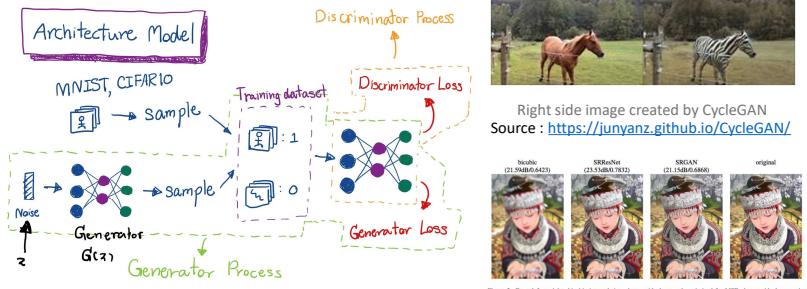
To create something from nothing is one of the greatest feelings, ... It's heaven.

By Prince

Python for Deep Learning

121

Generative Adversarial Networks : GANs



Python for Deep Learning

122

Generative Adversarial Networks : GANs



Picture: These people are not real – they were produced by our generator that allows control over different aspects of the image.

Source: <https://github.com/NVlabs/stylegan>

A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras (NVIDIA), Samuli Laine (NVIDIA), Timo Aila (NVIDIA)

<https://arxiv.org/abs/1812.04948>

Python for Deep Learning

123

Generative Adversarial Nets

Ian J. Goodfellow, Jean Pouget-Abadie*, Mehdi Mirza, Bing Xu, David Warde-Farley,

Sherjil Ozair†, Aaron Courville, Yoshua Bengio†

Département d'informatique et de recherche opérationnelle

Université de Montréal

Montréal, QC H3C 3J7

Abstract

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model G that captures the data distribution, and a discriminative model D that estimates

Paper : <https://arxiv.org/pdf/1406.2661.pdf>

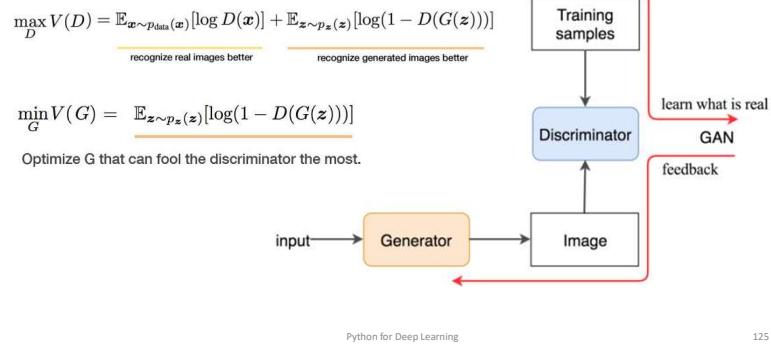
Python for Deep Learning

124

123

124

How GANs works?



125

Generative Adversarial Networks : GANs

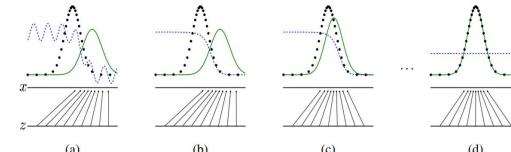


Figure 1: Generative adversarial nets are trained by simultaneously updating the discriminative distribution (D , blue, dashed line) so that it discriminates between samples from the data generating distribution (black, dotted line) p_{data} from those of the generative distribution p_g (G , green, solid line). The lower horizontal line is the domain from which \mathbf{z} is sampled, in this case uniformly. The horizontal line above is part of the domain of \mathbf{x} . The upward arrows show how the mapping $\mathbf{x} = G(\mathbf{z})$ imposes the non-uniform distribution p_g on transformed samples. G contracts in regions of high density and expands in regions of low density of p_g . (a) Consider an adversarial pair near convergence: p_g is similar to p_{data} and D is a partially accurate classifier. (b) In the inner loop of the algorithm D is trained to discriminate samples from data, converging to $D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$. (c) After an update to G , gradient of D has guided $G(\mathbf{z})$ to flow to regions that are more likely to be classified as data. (d) After several steps of training, if G and D have enough capacity, they will reach a point at which both cannot improve because $p_g = p_{\text{data}}$. The discriminator is unable to differentiate between the two distributions, i.e. $D(\mathbf{x}) = \frac{1}{2}$.

Python for Deep Learning

126

126

Generative Adversarial Networks : GANs

GANs is based on the [zero-sum non-cooperative game](#). In short, if one wins the other loses. A zero-sum game is also called [minimax](#). Your opponent wants to maximize its actions and your actions are to minimize them. In game theory, the GAN model converges when the discriminator and the generator reach a [Nash equilibrium](#). This is the optimal point for the minimax equation below.

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Zero-Sum Games

Zero-sum game: Rock-Paper-Scissor

P1 \ P2	Rock	Paper	Scissor
Rock	(0, 0)	(-1, 1)	(1, -1)
Paper	(1, -1)	(0, 0)	(-1, 1)
Scissor	(-1, 1)	(1, -1)	(0, 0)

Python for Deep Learning

127

127

What's Nash Equilibrium ?

In the game theory, the [Nash Equilibrium](#) is reached when no player will change its strategy after considering all possible strategy of opponents. i.e. in the Nash equilibrium, no one will change its decision even after we will all the player strategy to everyone. A game can have 0, 1 or multiple Nash Equilibria.

Non-zero-sum game: Prisoner's dilemma

		John	Stay Silent	Betray	
		Stay Silent	(-1, -1)	(-3, 0)	(-2, -2)
		Betray	(0, -3)	(-2, -2)	
					(-2, -2) is Nash equilibrium!

Python for Deep Learning

128

128

In other words, D and G play the following two-player minimax game with value function $V(G, D)$:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]. \quad (1)$$

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

```

for number of training iterations do
  for  $k$  steps do
    • Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
    • Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
    • Update the discriminator by ascending its stochastic gradient:
      
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(\mathbf{x}^{(i)}) + \log(1 - D(G(\mathbf{z}^{(i)})))]$$
.
  end for
  • Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
  • Update the generator by descending its stochastic gradient:
    
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(\mathbf{z}^{(i)})))$$
.
end for
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.
  
```

Python for Deep Learning

129

Generative Adversarial Networks : GANs

NIPS 2016 Tutorial:
Generative Adversarial Networks

Ian Goodfellow
OpenAI, ian@openai.com

Abstract

This report summarizes the tutorial presented by the author at NIPS 2016 on *generative adversarial networks* (GANs). The tutorial describes: (1) Why generative modeling is a topic worth studying, (2) how generative models work, and how GANs compare to other generative models, (3) the details of how GANs work, (4) research frontiers in GANs, and (5) state-of-the-art image models that combine GANs with other methods. Finally, the tutorial contains three exercises for readers to complete, and the solutions to these exercises.

Paper : <https://arxiv.org/pdf/1701.00160.pdf>

Python for Deep Learning

130

129

130

Training GANs is so HARD!

Many GAN models suffer the following major problems:

- **Non-convergence**: the model parameters oscillate, destabilize and never converge,
- **Mode collapse**: the generator collapses which produces limited varieties of samples,
- **Diminished gradient**: the discriminator gets too successful that the generator gradient vanishes and learns nothing,
- Unbalance between the generator and discriminator causing overfitting, &
- Highly sensitive to the hyperparameter selections.

Python for Deep Learning

131

Training GANs is so HARD!

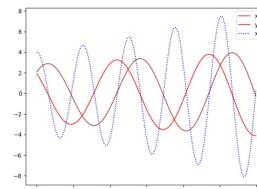
Many GAN models suffer the following major problems:

- **Non-convergence**: the model parameters oscillate, destabilize and never converge,

$$\min_B \max_A V(D, G) = \mathbf{x}\mathbf{y}$$

$$\Delta x = \alpha \frac{\partial(xy)}{\partial(x)}$$

$$\Delta y = -\alpha \frac{\partial(xy)}{\partial(y)}$$



Cost functions may not converge using gradient descent in a minimax game.

Python for Deep Learning

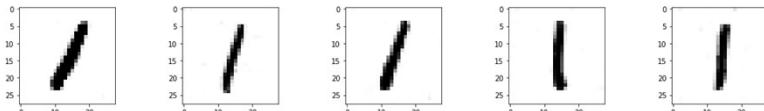
132

131

132

Training GANs is so HARD!

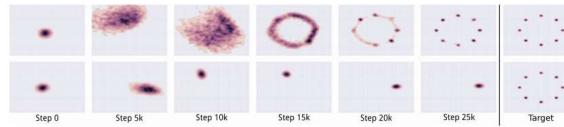
- **Mode collapse:** the generator collapses which produces limited varieties of samples,



Python for Deep Learning

133

Training GANs is so HARD!



Unrolled GANs are able to fit all of the modes of a mixture of Gaussians in a two-dimensional space.
Paper Unrolled GANs : <https://arxiv.org/pdf/1611.02163.pdf>



Python for Deep Learning

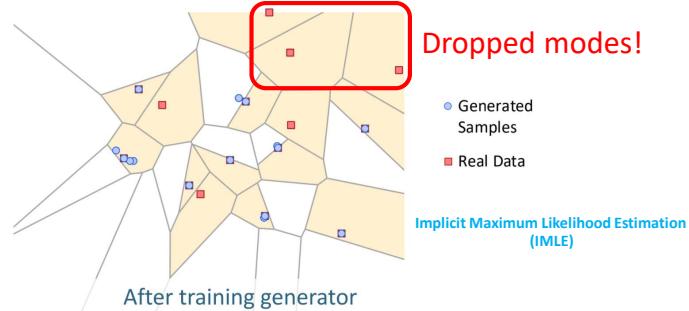
134

133

134

Why Mode Collapse Happens

Source : https://drive.google.com/file/d/1PV4YN3OQprww4BCDwB9XWMUIz_mbdDab/view



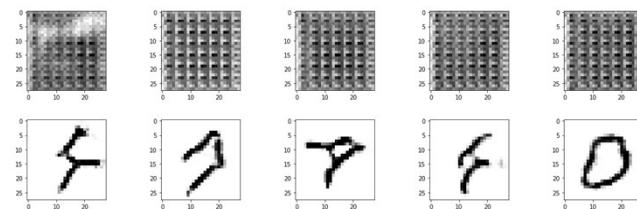
Python for Deep Learning

135

135

Training GANs is so HARD!

- **Diminished gradient:** the discriminator gets too successful that the generator gradient vanishes and learns nothing,



Python for Deep Learning

136

136

In other words, D and G play the following two-player minimax game with value function $V(G, D)$:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]. \quad (1)$$

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log(1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(\mathbf{z}^{(i)}))). \quad \nabla_{\mathbf{w}_G} \frac{1}{n} \sum_{i=1}^n \log(D(G(\mathbf{z}^{(i)})))$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Python for Deep Learning

137

Training GANs is so HARD!

Many GAN models suffer the following major problems:

- **Non-convergence**: the model parameters oscillate, destabilize and never converge,
- **Mode collapse**: the generator collapses which produces limited varieties of samples,
- **Diminished gradient**: the discriminator gets too successful that the generator gradient vanishes and learns nothing,
- Unbalance between the generator and discriminator causing overfitting, &
- Highly sensitive to the hyperparameter selections.

Python for Deep Learning

138

Some improvements of GANs

• Improving the Network design

- DCGAN (Deep convolutional generative adversarial networks), **pop!**
- CGAN, InfoGAN (using labels to improve GAN),
- SGAN (Stacked GAN),
- StyleGAN, etc.

Python for Deep Learning

139

Some improvements of GANs

• Improving the cost function

Finally, we did not find evidence that any of the tested algorithms consistently outperforms the original one.

A 2017 Google Brain paper "Are GANs Created Equal?"
Paper : <https://arxiv.org/pdf/1711.10337.pdf>

Name	Value Function
GAN	$L_D^{\text{GAN}} = \mathbb{E}[\log(D(x))] + \mathbb{E}[\log(1 - D(G(z)))]$ $L_G^{\text{GAN}} = \mathbb{E}[\log(D(G(z)))]$
LSGAN	$L_D^{\text{LSGAN}} = \mathbb{E}[(D(x) - 1)^2] + \mathbb{E}[D(G(z))^2]$ $L_G^{\text{LSGAN}} = \mathbb{E}[(D(G(z)) - 1)^2]$
WGAN	$L_D^{\text{WGAN}} = \mathbb{E}[D(x)] - \mathbb{E}[D(G(z))]$ $L_G^{\text{WGAN}} = \mathbb{E}[D(G(z))]$ $W_p \leftarrow \text{clip_by_value}(W_p, -0.01, 0.01)$
WGANGP	$L_D^{\text{WGANGP}} = L_D^{\text{WGAN}} + \lambda E[(D(ax - (1 - \alpha G(x))) - 1 ^2)]$ $L_G^{\text{WGANGP}} = L_G^{\text{WGAN}}$
DRAGAN	$L_D^{\text{DRAGAN}} = L_D^{\text{GAN}} + \lambda E[(D(ax - (1 - \alpha x_p)) - 1 ^2)]$ $L_G^{\text{DRAGAN}} = L_G^{\text{WGAN}}$
COGAN	$L_D^{\text{COGAN}} = \mathbb{E}[\log(D(x, c))] + \mathbb{E}[\log(1 - D(G(z), c))]$ $L_G^{\text{COGAN}} = \mathbb{E}[\log(D(G(z), c))]$
infoGAN	$L_D^{\text{infoGAN}} = L_D^{\text{GAN}} - \mathcal{L}_{KL}(G, c)$ $L_G^{\text{infoGAN}} = L_G^{\text{GAN}} - \mathcal{L}_{KL}(G, c)$
ACGAN	$L_D^{\text{ACGAN}} = L_D^{\text{GAN}} + E[P(\text{class} = c x)] + E[P(\text{class} = c G(x))]$ $L_G^{\text{ACGAN}} = L_G^{\text{GAN}} + E[P(\text{class} = c G(x))]$
EBOGAN	$L_D^{\text{EBOGAN}} = D_{\text{KL}}(x + \max(0, m - D_{\text{KL}}(G(x)))$ $L_G^{\text{EBOGAN}} = D_{\text{KL}}(G(x) + \lambda \cdot PT)$
BEGAN	$L_D^{\text{BEGAN}} = D_{\text{KL}}(x - k_i D_{\text{KL}}(G(x))$ $L_G^{\text{BEGAN}} = D_{\text{KL}}(G(x) - k_i + \lambda(y D_{\text{KL}}(x) - D_{\text{KL}}(G(x))))$

Python for Deep Learning

140

139

140

Improved Techniques for Training GANs

Tim Salimans Ian Goodfellow Wojciech Zaremba Vicki Cheung
 tim@openai.com ian@openai.com woj@openai.com vicki@openai.com

Alec Radford Xi Chen
 alec@openai.com peter@openai.com

Abstract

We present a variety of new architectural features and training procedures that we apply to the generative adversarial networks (GANs) framework. Using our new techniques, we achieve state-of-the-art results in semi-supervised classification on MNIST, CIFAR-10 and SVHN. The generated images are of high quality as con-

Paper : <https://arxiv.org/pdf/1606.03498.pdf>

Python for Deep Learning 141

141

Hi limpapat,

For centuries, we've relied on experts to spot art forgeries. Using GANs, algorithms can separate the fakes from the fakes by learning to imitate an artist's style.

In our newest getting started competition, you will train GANs to mimic the works of Claude Monet. Follow the style of an existing photo or create new Monets from scratch.

This Getting Started competition will be ongoing at any time. But [TPU Star Awards](#) for extra TPU credits until October 31, 2020!

Put some art into (data) science, improve your knowledge classifiers into believing you've created a true Monet.

[Join the Competition](#)

Python for Deep Learning 142

142

Generative Adversarial Networks : GANs



BLOG
NVIDIA GAUGAN

Source : <https://www.youtube.com/watch?v=OGGjXG562WU>

The GANs Zoo : <https://github.com/hindupuravinash/the-gan-zoo>

GANs Applications : <https://github.com/nashory/gans-awesome-applications>

Python for Deep Learning 143

143