# Introduction to Deep Learning

Some slides and images are taken from:

David Wolfe Corne Wikipedia Geoffrey A. Hinton https://www.macs.hw.ac.uk/~dwcorne/Teaching/introdl.ppt

### Feedforward networks for function approximation and classification



- Feedforward networks
- Deep tensor networks
- ConvNets



### A single artificial neuron



### A single artificial neuron



### Activation functions

https://en.wikipedia.org/wiki/Activation\_function



### Activation functions

https://en.wikipedia.org/wiki/Activation\_function



### Activation functions

https://en.wikipedia.org/wiki/Activation\_function



Currently most widely used. Empirically easier to train and results in sparse networks.

Nair and Hinton: Rectified linear units improve restricted Boltzmann machines ICLM'10, 807-814 (2010) Glorot, Bordes and Bengio: Deep sparse rectifier neural networks. PMLR 15:315-323 (2011)

### Perceptron (Rosenblatt 1957)

Used as a binary classifier - equivalent to support vector machine (SVM)



Train weights using examples  $D = \{(\mathbf{x}_1, d_1), ..., (\mathbf{x}_S, d_S)\}$ :

- 1. Initialize  $\mathbf{w}$  with random values
- 2. For each example  $(\mathbf{x}_j, d_j)$ :
  - (a) Calculate output:  $y_j = f(\mathbf{w}^\top \mathbf{x}_j)$
  - (b) Update weights:  $\mathbf{w} \leftarrow \mathbf{w} + (d_j y_j)\mathbf{x}_j$
- 3. Repeat 2 until convergence.

#### Perceptron (Rosenblatt 1957)





Slide credit : Geoffrey Hinton

Compare outputs with **Back-propagate** correct answer to get error signal to get error signal derivatives for learning outputs hidden layers input vector

Slide credit : Geoffrey Hinton

Error at output for a given example:

$$E = \frac{1}{2} \sum_{j} (y_j - d_j)^2$$



Slide credit : Geoffrey Hinton

Error at output for a given example:

$$E = \frac{1}{2} \sum_{j} (y_j - d_j)^2$$

Error sensitivity at output neuron j:

$$\frac{\partial E}{\partial y_j} = y_j - d_j$$

Backpropagate error sensitivity to neuron i:

$$\frac{\partial E}{\partial y_i} = \sum_j \frac{\partial E}{\partial y_j} f'(x_j) w_{ji}$$

Sensitivity on weight w<sub>ji</sub>:

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial x_j} y_i$$



Slide credit : Geoffrey Hinton

Error at output for a given example:

$$E = \frac{1}{2} \sum_{j} (y_j - d_j)^2$$

Error sensitivity at output neuron j:

$$\frac{\partial E}{\partial y_j} = y_j - d_j$$

Backpropagate error sensitivity to neuron i:

$$\frac{\partial E}{\partial y_i} = \sum_j \frac{\partial E}{\partial y_j} f'(x_j) w_{ji}$$

Sensitivity on weight w<sub>ji</sub>:

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial x_j} y_i$$

Update weight w<sub>ji</sub>:

$$\Delta w_{ji} = -\epsilon \frac{\partial E}{\partial w_{ji}}$$
learning rate



Initial random weights













### Nonlinear versus linear models

NNs use nonlinear f(x) so they can draw complex boundaries, but keep the data unchanged Kernel methods only draw straight lines, but transform the data first in a way that makes it linearly separable

### **Universal Representation Theorem**

- Networks with a single hidden layer can represent any function F(x) with arbitrary precision in the large hidden layer size limit
- However, that doesn't mean, networks with single hidden layers are efficient in representing arbitrary functions. For many datasets, deep networks can represent the function F(x) even with narrow layers.

### What does a neural network learn?





Figure 1.2: Examples of handwritten digits from postal envelopes.



## Feature detectors



### What is this unit doing?



Figure 1.2: Examples of handwritten digits from postal envelopes.





Hidden layer units become self-organized feature detectors



![](_page_26_Figure_1.jpeg)

![](_page_27_Figure_1.jpeg)

![](_page_28_Figure_1.jpeg)

![](_page_29_Figure_1.jpeg)

![](_page_30_Picture_0.jpeg)

Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.

What features might you expect a good NN to learn, when trained with data like this?

![](_page_31_Figure_0.jpeg)

Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.

![](_page_32_Picture_0.jpeg)

Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.

![](_page_33_Figure_0.jpeg)

Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.

![](_page_34_Figure_0.jpeg)

But what about position invariance? Our example unit detectors were tied to specific parts of the image

## **Deep Networks**

![](_page_35_Picture_1.jpeg)

### successive layers can detect higher-level features

![](_page_36_Picture_1.jpeg)

### successive layers can detect higher-level features

![](_page_37_Picture_1.jpeg)

### So: multiple layers make sense

![](_page_38_Figure_1.jpeg)

### So: multiple layers make sense

#### Multiple layers are also found in the brain, e.g. visual cortex

![](_page_39_Figure_2.jpeg)

![](_page_39_Figure_3.jpeg)

But: until recently deep networks could not be efficiently trained

![](_page_40_Picture_1.jpeg)

# 2006: The Deep Breakthrough

![](_page_41_Picture_1.jpeg)

- Hinton, Osindero & Teh
   « <u>A Fast Learning</u>
   <u>Algorithm for Deep</u>
   <u>Belief Nets</u> », Neural
   Computation, 2006
- Bengio, Lamblin, Popovici, Larochelle « <u>Greedy Layer-Wise</u> <u>Training of Deep</u> <u>Networks</u> », *NIPS'2006*
  - Ranzato, Poultney, Chopra, LeCun « Efficient Learning of Sparse Representations with an Energy-Based Model », NIPS'2006

## **Convolutional Neural Networks**

![](_page_42_Figure_1.jpeg)

Compared to standard feedforward neural networks with similarly-sized layers,

- CNNs have much fewer connections and parameters
- and so they are easier to train,
- while their theoretically-best performance is likely to be only slightly worse.

#### LeNet 5

Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998 Convolutional Kernel / Filter

0	1	2
2	2	0
0	1	2

#### Apply convolutions

![](_page_44_Figure_3.jpeg)

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	30	$2_1$	$1_2$	0
0	02	$1_2$	30	1
3	1,0	$2_1$	$2_2$	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
00	0,	$1_2$	3	1
32	$1_2$	$2_0$	2	3
$2_0$	0,	02	2	2
2	0	0	0	1

12.0	12.0	17.0
12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	00	1,	32	1
3	$1_2$	$2_2$	$2_0$	3
2	00	0,	$2_2$	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

![](_page_45_Figure_0.jpeg)

- Input: 32x32 pixel image. Largest character is 20x20 (All important info should be in the center of the receptive field of the highest level feature detectors)
- Cx: Convolutional layer
- Sx: Subsample layer
- Fx: Fully connected layer
- Black and White pixel values are normalized:
   E.g. White = -0.1, Black =1.175 (Mean of pixels = 0, Std of pixels =1)

## Convolutional filters perform image processing

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

## Example: filters in face recognition

![](_page_47_Picture_1.jpeg)

![](_page_47_Picture_2.jpeg)

![](_page_47_Picture_3.jpeg)

First Layer Representation

Second Layer Representation Third Layer Representation

## **MNIST** dataset

540,000 artificial distortions

+ 60,000 original

Test error: 0.8%

60,000 original datasets Test error: 0.95%

![](_page_48_Figure_6.jpeg)

## **Misclassified examples**

![](_page_49_Picture_1.jpeg)

## Applications to molecular systems

![](_page_50_Picture_1.jpeg)

### 1) Learning to represent (effective) energy function

#### Behler-Parrinello network

![](_page_51_Figure_2.jpeg)

Behler and Parrinello, PRL 98, 146401 (2007)

### 1) Learning to represent (effective) energy function

![](_page_52_Figure_1.jpeg)

![](_page_52_Figure_2.jpeg)

Schuett, Arbabzadah, Chmiela, Müller & Tkatchenko, Nature Communications 8, 13890 (2017)

### 2) Generator networks

![](_page_53_Figure_1.jpeg)

Gómez-Bombarelli, ..., Aspuru-Guzik: Automatic Chemical Design using Variational Autoencoders (2016)

### 2) Generator networks

![](_page_54_Picture_1.jpeg)

Figure 3: a). Random sampling. Molecules decoded from randomly-sampled points in the latent space of a variational autoencoder, near to a given molecule (aspirin [2-(acetyloxy)benzoic acid], highlighted in blue). b). Interpolation. Two-dimensional interpolation between four random points in in drug-like VAE. Decodings of interpolating linearly between the latent representations of the four molecules in the corners.

Gómez-Bombarelli, ..., Aspuru-Guzik: Automatic Chemical Design using Variational Autoencoders (2016)

### 2) Generator networks

![](_page_55_Figure_1.jpeg)

Gómez-Bombarelli, ..., Aspuru-Guzik: Automatic Chemical Design using Variational Autoencoders (2016)

### 3) VAMPnets

![](_page_56_Figure_1.jpeg)

Mardt, Pasquali, Wu & Noé: VAMPnets - deep learning of molecular kinetics (2017)

### 3) VAMPnets

![](_page_57_Figure_1.jpeg)

1	1.>1	1->2	1->3	1->4	1->5	1->6
0	Parented	,				
1	2->1	2->2	2->3	2->4	2->5	2->6
0	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	L		,		
1	3->1	3->2	3->3	3->4	3->5	3->6
0	,		L	,		
1	4->1	4->2	4->3	4->4	4->5	4->6
0	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	,	,	<b>h</b>		
1	5->1	5.>2	5.>3	5->4	4444A	5->6
0				-+++++++		
1	6->1	6->2	6.>3	6->4	, <b>∤#</b> #\$\$	6->6
0				-++++++++	/	
	0.1 0.3	0.1 0.3	0.1 0.3	0.1 0.3	0.1 0.3	0.1 0.3

Mardt, Pasquali, Wu & Noé: VAMPnets - deep learning of molecular kinetics (2017)