Machine learning for hearing research

Kang-Hun Ahn

Department of Physics Chungnam National University

I. Introduction

- II. Hearing in noise: machine learning simulation
- III. Speech recognition faster than FFT using machine learning
- IV. Tinnitus as a neural net property
- V. Unsupervised learning: Dictionary learning of sound



What Is Deep Learning?



Deep learning

From Wikipedia, the free encyclopedia

Deep learning (*deep machine learning*, or *deep structured learning*, or *hierarchical learning*, or sometimes *DL*) is a branch of <u>machine learning</u> based on a set of <u>algorithms</u> that attempt to model high-level abstractions in data by using model architectures, with complex structures or otherwise, composed of multiple <u>non-linear transformations</u>.^{[1](p198)[2][3][4]}

REVIEW

doi:10.1038/nature14539

Deep learning

Yann LeCun^{1,2}, Yoshua Bengio³ & Geoffrey Hinton^{4,5}

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other domains such as drug discovery and genomics. Deep learning discovers intricate structure in large data sets by using the backpropagation algorithm to indicate how a machine should change its internal parameters that are used to compute the representation in each layer from the representation in the previous layer. Deep convolutional nets have brought about breakthroughs in processing images, video, speech and audio, whereas recurrent nets have shone light on sequential data such as text and speech.

436 | NATURE | VOL 521 | 28 MAY 2015





Biological neuron architecture



From neuron to brain



The information flow is one-way (directional).



Artificial Neuron

Perceptron



Frank Rosenblatt (1927-1971)





Problem: a small change in the weights or bias of any single perceptron can cause sometimes cause the output of that perceptron to completely flip, say 0 to 1.



Sigmoid neuron



Types of Learning

Supervised (inductive) learning

· Training data includes desired outputs

Unsupervised learning

· Training data does not include desired outputs

Semi-supervised learning

· Training data includes a few desired outputs

The architecture of neural network



Role of hidden layer





XOR problem : Monolayer is not enough. Hidden layer is required.



Training : Minimizing an object function (error)

$$C(w,b)\equiv \sum_{x} \left\|y(x)-a
ight\|^{2}$$

Desired output

Calculated output





$$abla C \equiv \left(rac{\partial C}{\partial v_1}\,,\ldots,rac{\partial C}{\partial v_m}
ight)^T$$

$$\Delta v = \left(\Delta v_1, \dots, \Delta v_m
ight)^T$$

$$\Delta v = -\eta
abla C$$

Learning : Repeated application of

$$egin{aligned} w_k \ o \ w'_k &= w_k - \eta \, rac{\partial C}{\partial w_k} \ b_l \ o \ b'_l &= b_l - \eta \, rac{\partial C}{\partial b_l} \,. \end{aligned}$$

Implementation:

- 1. How to calculate the gradients ?
- 2. Batch learning, online learning ?
- **3.** Learning rate η ?
- 4. Overtraining.

The error of neuron δ_j^l



$$\delta^l_j \equiv rac{\partial C}{\partial z^l_j} \qquad z^l_j = \sum_k w^l_{jk} a^{l-1}_k + b^l_j$$

,where the error cost function is

$$C = rac{1}{2} \, \|y-a^L\|^2 = rac{1}{2} \sum_j (y_j-a_j^L)^2 \, ,$$

The error of output layer neuron is clearly

$$\delta^L_j = rac{\partial C}{\partial a^L_j} \, \sigma'(z^L_j)$$

Back propagation

A simple linear equation for the error in terms of the error of next layer

$$\begin{split} \delta_j^l &= \sum_k w_{kj}^{l+1} \delta_k^{l+1} \sigma'(z_j^l) \\ \text{Proof}) \quad \delta_j^l &= \frac{\partial C}{\partial z_j^l} \\ &= \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} \\ &= \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1}, \\ \text{where} \quad z_k^{l+1} &= \sum_j w_{kj}^{l+1} a_j^l + b_k^{l+1} = \sum_j w_{kj}^{l+1} \sigma(z_j^l) + b_k^{l+1}, \\ &\qquad \frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \sigma'(z_j^l) \end{split}$$

1. How to calculate the gradients? Solution:

The equations of backpropagation
1.
$$\delta^{L} = \nabla_{a}C \odot \sigma'(z^{L})$$

2. $\delta^{l} = ((w^{l+1})^{T}\delta^{l+1}) \odot \sigma'(z^{l})$
3. $\frac{\partial C}{\partial b_{j}^{l}} = \delta_{j}^{l}$
4. $\frac{\partial C}{\partial w_{jk}^{l}} = a_{k}^{l-1}\delta_{j}^{l}$

Hardamard product

- Input x: Set the corresponding activation a¹ for the input layer.
- 2. Feedforward: For each $l = 2, 3, \ldots, L$ compute

 $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$.

- 3. **Output error** δ^L : Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$.
- 4. Backpropagate the error: For each l = L 1, L 2, ..., 2compute $\delta^{l} = ((w^{l+1})^{T} \delta^{l+1}) \odot \sigma'(z^{l})$.
- 5. **Output:** The gradient of the cost function is given by $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j^l} = \delta_j^l.$
- 6. **Gradient descent:** For each l = L, L 1, ..., 2 update the weights according to the rule $w^l \to w^l \frac{\eta}{m} \sum_x \delta^{x,l} (a^{x,l-1})^T$, and the biases according to the rule $b^l \to b^l \frac{\eta}{m} \sum_x \delta^{x,l}$.

2. Batch learning or online learning?

Suppose you have different data sets A,B, C

• Batched Learning
$$\delta \omega = -\eta \sum_{\alpha = A, B, C} \left\langle \alpha \left| \frac{\partial C(\omega)}{\partial \omega} \right| \alpha \right\rangle$$
$$\omega \xrightarrow{A, B, C} \omega + \delta \omega$$

· Online Learning



2. Batched learning or online learning?

- These two processes are mathematically different.
- Both of them work well for training.



Too small learning rate gives too slow convergence.

Too large learning rate does not give convergence.

There exists an range of optimal learning rate which must be found through experience.

4. Overtraining

Too much training will hurt you ~~



Most people in this field know this phenomenon but there is no mathematical proof.



Recipe for (Supervised) Deep Research Learning with Big Data

Microsoft[®]



Technical problems

Why gradients tend to vanish for DNN

· To illustrate to problem, let's use matrix form of error BP:

$$\begin{split} \delta^{l} &= ((w^{l+1})^{T} \delta^{l+1}) \odot \sigma'(z^{l}) \\ & \bullet \\ &= ((w^{\prime+2})^{T} \left[((w^{\prime+2})^{T} \delta^{\prime+2}) \odot \sigma'(z^{l})^{\prime} \right]) \odot \sigma'(z^{l})^{\prime} \end{split}$$

This problem becomes more serious as the layer becomes deeper.

- · So even if forward pass in nonlinear, the error backprop is a linear process
- It suffers from all problems associated with linear processes
- Many terms of σ (1- σ)
- · In addition, many terms in the product of W's
- · If any sigmoid unit saturates in either direction, the error gradient becomes zero
- If ||W||<1, the product will shrink fast for high depths
- If ||W||>1, the product may grow fast for high depths

From Le Deng's lecture note

1.
$$\delta^{L} = \nabla_{a}C \odot \sigma'(z^{L})$$

2. $\delta^{l} = ((w^{l+1})^{T} \delta^{l+1}) \odot \sigma'(z^{l})$
3. $\frac{\partial C}{\partial b_{j}^{l}} = \delta_{j}^{l}$
4. $\frac{\partial C}{\partial w_{jk}^{l}} = a_{k}^{l-1} \delta_{j}^{l}$

Smallness of σ' slow down the learning process

variations

Log-likehood cost function

$$C = -\ln a_y^L$$

This is useful to avoid the slow-down problem as

T

$$\begin{aligned} \frac{\partial C}{\partial w_{jk}^L} &= a_k^{L-1} (a_j^L - y_j) \\ \frac{\partial C}{\partial b_j^L} &= a_j^L - y_j \end{aligned}$$

Softmax outlayer

$$a_j^L = \frac{e^{z_j^L}}{\sum_k e^{z_k^L}}$$

Can be interpreted as probability of the output j as sum of a_i^L is equal to 1.



-mathematically they implement static input-output mapping

Multi-layer perception(MLP) can approximate arbitrary nonlinear functions with arbitrary precision

-Most popular supervised training algorithm: backpropagation algorithm

-Most (95%?) of neural network publication concern feedforward net

-have proven useful in many practical applications as pattern classifications

Recurrent network

- all biological neural networks are recurrent
- mathematically, RNNs implement dynamical systems
- basic theoretical result: RNNs can approximate arbitrary (term needs some qualification) dynamical systems with arbitrary precision ("universal approximation property")
- · several types of training algorithms are known, no clear winner
- theoretical and practical difficulties by and large have prevented practical applications so far



Formal description of RNN





Updates of internal units

 $\mathbf{x}(n+1) = \mathbf{f}(\mathbf{W}^{in}\mathbf{u}(n+1) + \mathbf{W}\mathbf{x}(n) + \mathbf{W}^{back}\mathbf{y}(n))$ internal Input output

The output

 $\mathbf{y}(n+1) = \mathbf{f}^{out}(\mathbf{W}^{out}(\mathbf{u}(n+1),\mathbf{x}(n+1)))$

f= tanh or 1

Training of recurrent network

Backpropagation through time (BPTT) method





Unfold the recurrent network in time, by stacking identical copies of RNN !

Teacher data $\mathbf{u}(n) = (u_1(n), ..., u_K(n))'$, $\mathbf{d}(n) = (d_1(n), ..., d_L(n))'$ n = 1, ...T

The error to be minimized $E = \sum_{n=1,\dots,T} \|\mathbf{d}(n) - \mathbf{y}(n)\|^2 = \sum_{n=1,\dots,T} E(n)$ where $\mathbf{y}(n+1) = \mathbf{f}^{out} (\mathbf{W}^{out} (\mathbf{u}(n+1), \mathbf{x}(n+1))_{\star})$

 $\mathbf{x}(n+1) = \mathbf{f}(\mathbf{W}^{in}\mathbf{u}(n+1) + \mathbf{W}\mathbf{x}(n) + \mathbf{W}^{back}\mathbf{y}(n)),$

Then the algorithm is straightforward feedforward backpropagation algorithm.

$$\delta_{j}(T) = (d_{j}(T) - y_{j}(T)) \frac{\partial f(u)}{\partial u}\Big|_{u=z_{j}(T)}$$

error for the output units

$$\delta_{i}(T) = \left[\sum_{j=1}^{L} \delta_{j}(T) w_{ji}^{out}\right] \frac{\partial f(u)}{\partial u}\Big|_{u=z_{i}(n)}$$

error for internal units x_i (t)

$$\delta_{j}(n) = \left[(d_{j}(n) - y_{j}(n)) + \sum_{i=1}^{N} \delta_{i}(n+1) w_{ij}^{back} \right] \frac{\partial f(u)}{\partial u} \Big|_{u=z_{j}(n)} \text{ error for the output units of earlier layer}$$

$$\delta_{i}(n) = \left[\sum_{j=1}^{N} \delta_{j}(n+1)w_{ji} + \sum_{j=1}^{L} \delta_{j}(n)w_{ji}^{out}\right] \frac{\partial f(u)}{\partial u}\Big|_{u=z_{i}(n)} \quad \text{error for the internal units of earlier times}$$

Adjustment

$$\begin{split} new \ w_{ij} &= w_{ij} + \gamma \sum_{n=1}^{T} \delta_i(n) x_j(n-1) \quad [use \ x_j(n-1) = 0 \ for \ n = 1] \\ new \ w_{ij}^{in} &= w_{ij}^{in} + \gamma \sum_{n=1}^{T} \delta_i(n) u_j(n) \\ new \ w_{ij}^{out} &= w_{ij}^{out} + \gamma \times \begin{cases} \sum_{n=1}^{T} \delta_i(n) u_j(n), & \text{if } j \ refers \ to \ input \ unit \\ \sum_{n=1}^{T} \delta_i(n) x_j(n), & \text{if } j \ refers \ to \ hidden \ unit \\ new \ w_{ij}^{back} &= w_{ij}^{back} + \gamma \sum_{n=1}^{T} \delta_i(n) y_j(n-1) \quad [use \ y_j(n-1) = 0 \ for \ n = 1] \end{cases}$$

I. Introduction

II. Hearing in noise: machine learning simulation

- III. Speech recognition faster than FFT using machine learning
- IV. Tinnitus as a neural net property
- V. Unsupervised learning: Dictionary learning of sound

How humans can understand speech in spite of background noise is a key issue in auditory science



Automatic speech (complex sound) recognition is still far from human speech recognition

cocktail party problem

What is the effect of noise on the hearing of complex sounds such as syllables?







150 syllables with high probabilities and 150 syllables with low probabilitiesUsual syllableUnusual syllable

1st Clinical trial

- 23 men and 26 women
- Male age range is 20~26, female age range is 19~24
- 300 Korean syllables. (60 dB(a))
- 5 white noise levels (38 dB(a), 42 dB(a), 47 dB(a), 52

dB(a), 57 dB(a))



The behavior looks like stochastic resonance, but

Stochastic resonance

In clinical trial, noise-enhanced syllable hearing appeared only for the specific syllables.

Noise-enhanced hearing syllables are rarely used syllables



Syllable (in increasing correctness order)

However, these are not the cause of this phenomenon

Noise can improve other cognitive abilities

Tactile sensation D/A Converter Potentiomete Electrodes Stimulus Isolator mes-Weinstein Presentation 0.6 Without Electrical Noise With Electrical Noise 0.5 04 0.3 õ 02 0.1 Subj I Subj 2 Subj 3 Subj 4 Subj 5 Subj 6 Subj 7 Subj 8 Subj 9 ** ** * ** ** * p < 0.05 ** ▷<0.01 Dhruv, Neel T., et al. "Enhancing tactile sensation in older

adults with electrical noise stimulation." Neuroreport 13.5

(2002): 597-600.

Stimulus Force Plaform

Sway measures	No stimulation	Stimulation	Improvement (%)
ML SD (mm)	4.6 ± 0.7	4.5 ± 0.7	3.8
ML max (mm)	13.0 + 2.2	13.1 + 2.6	- 1.1
AP SD (mm)	6.3 + 1.3	6.1 + 1.4	3.7
AP max (mm)	20.0 + 5.6	19.0 + 5.3	5.4
Mean radius (mm)	6.9 + I.I	6.7 + 1.2	3.3
Path length (mm/s)	37.1 ± 5.0	35.9 ± 5.8	3.1
Swept area (mm ² /s)	88.8 ± 25.5	81.9 ± 24.7	7.8

Gravelle, Denise C., et al. "Noise-enhanced balance control in older adults."Neuroreport 13.15 (2002): 1853-1856.

Vision





Moss, Frank, Lawrence M. Ward, and Walter G. Sannita. "Stochastic resonance and sensory information processing: a tutorial and review of application." Clinical neurophysiology 115.2 (2004): 267-281.

Addition of an appropriate amount of noise can actually improve cognitive abilities

Balance control

Appropriate amount of noise can help hearing pure tone sounds by lowering threshold



F. Zeng, Q. Fu, R. Morse, Brain Res. 869 (2000) 251.

Appropriate amount of noise can help hearing pure tone sounds by lowering threshold





Long, Zhang-Cai, et al. Physics Letters A323.5 (2004): 434-438.

Correct percentages of detecting pure tones of different amplitudes with additive noise



What happened?

Response bias ? Individual disorder ? More attention with noise?

2nd Clinical trial

- 15 men and 22 women
- 200 Korean syllables with 4 levels of white noise (syllable 201~300 from the 1st clinical trial and the 57 dB(a))
- The order of the syllable was different for each subject (Remove response bias)



2nd Clinical trial



1st and 2nd clinical trial syllables have a similar increasing structure.

Syllable dependence or individual dependence

If the effect is due to individual disorders, there must b strong overlap between the people who have this effect for different syllables.

Number of overlap k between two randomly chosen groups of number of people m and n among a total N





This overlap number is small and it can happen by chance

The application of the machine learning to speech recognition



Syllable recognition test in automatic speech recognition

If noise-induced hearing enhancement of syllables comes from the **central nervous system**, it is natural to expect this phenomenon arise also in an **artificial speech recognition systems**.



The noise-induced hearing enhancement also found in TTS-STT system

Syllable recognition test in artificial speech recognition



Can we see the same phenomenon in **feedforward neural network**?



Training : Minimizing an object function (error)

$$C(w,b)\equiv \sum_{x} \left\|y(x)-a
ight\|^{2}$$

Desired output

Calculated output





Back propgation

The equations of backpropagation 1. $\delta^L = \nabla_a C \odot \sigma'(z^L)$ 2. $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ 3. $\frac{\partial C}{\partial b_j^l} = \delta_j^l$ 4. $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$

Hardamard product

- Input x: Set the corresponding activation a¹ for the input layer.
- 2. Feedforward: For each $l = 2, 3, \ldots, L$ compute

 $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$.

- 3. **Output error** δ^L : Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$.
- 4. Backpropagate the error: For each l = L 1, L 2, ..., 2compute $\delta^{l} = ((w^{l+1})^{T} \delta^{l+1}) \odot \sigma'(z^{l})$.
- 5. **Output:** The gradient of the cost function is given by $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j^l} = \delta_j^l.$
- 6. **Gradient descent:** For each l = L, L 1, ..., 2 update the weights according to the rule $w^l \to w^l \frac{\eta}{m} \sum_x \delta^{x,l} (a^{x,l-1})^T$, and the biases according to the rule $b^l \to b^l \frac{\eta}{m} \sum_x \delta^{x,l}$.

Undertraining vs Overtraining



Most people in this field know this phenomenon but there is no mathematical proof.

- Mel scale

MFCC: Mel scale and Mel Frequency Cepstral Coefficient

Humans are much better at discerning small changes in pitch at low frequencies than they are at high frequencies.

$$m(f) = 1125 \ln\left(1 + \frac{f}{700}\right) = 2595 \log_{10}\left(1 + \frac{f}{700}\right)$$

- MFCC

speech signal



1. Take the Discrete Fourier Transform of each frame.

$$S_i(f) = \sum_{n=1}^{N} s_i(n)h(n)e^{-\frac{i2\pi fn}{N}} \quad (1 \le f \le F)$$

where $S_i(n)$ is the *i*th time – domain frame and h(n) is Hamming window function.

2. Multiply each filterbank with the power spectrum and then add up the coefficients.

$$\tilde{S}_{ij} = \sum_{f=0}^{N/2} S_i(f) M_j(f)$$

where $M_i(f)$ is the jth mel filterbank.

3. Take the log of each \tilde{S}_{ij} and take the Discrete Cosine Transform.

$$c_i(n) = \frac{1}{\sqrt{J}} \sum_{j=0}^{J} \log(\tilde{S}_{ij}) \cos(\frac{(2j+1)n}{2J})$$

where J is the number of mel filterbanks.



MFCC

The pattern for different phonemes are more *visible* in MFCC.



Noise-induced hearing enhancement appears in insufficiently trained network



Syllable *o* have a maximum correctness at a specific noise intensity in the feedforward neural network.



Mathematical model made not that easy

$$S_{i}(f) = \sum_{n=1}^{N} s_{i}(n)h(n)e^{-\frac{i2\pi fn}{N}} \quad (1 \le f \le F)$$

$$\tilde{S}_{ij} = \sum_{f=0}^{N/2} S_{i}(f)M_{j}(f)$$

$$c_{i}(n) = \frac{1}{\sqrt{J}} \sum_{j=0}^{J} \log(\tilde{S}_{ij})\cos(\frac{(2j+1)n}{2J})$$

$$C(w,b) \equiv \frac{1}{2n} \sum_{x} ||y(x) - a||^{2}$$

$$w_{k} \rightarrow w_{k}' = w_{k} - \eta \frac{\partial C}{\partial w_{k}}$$

$$b_{l} \rightarrow b_{l}' = b_{l} - \eta \frac{\partial C}{\partial b_{l}}.$$