

Backpropagation, Neural Network

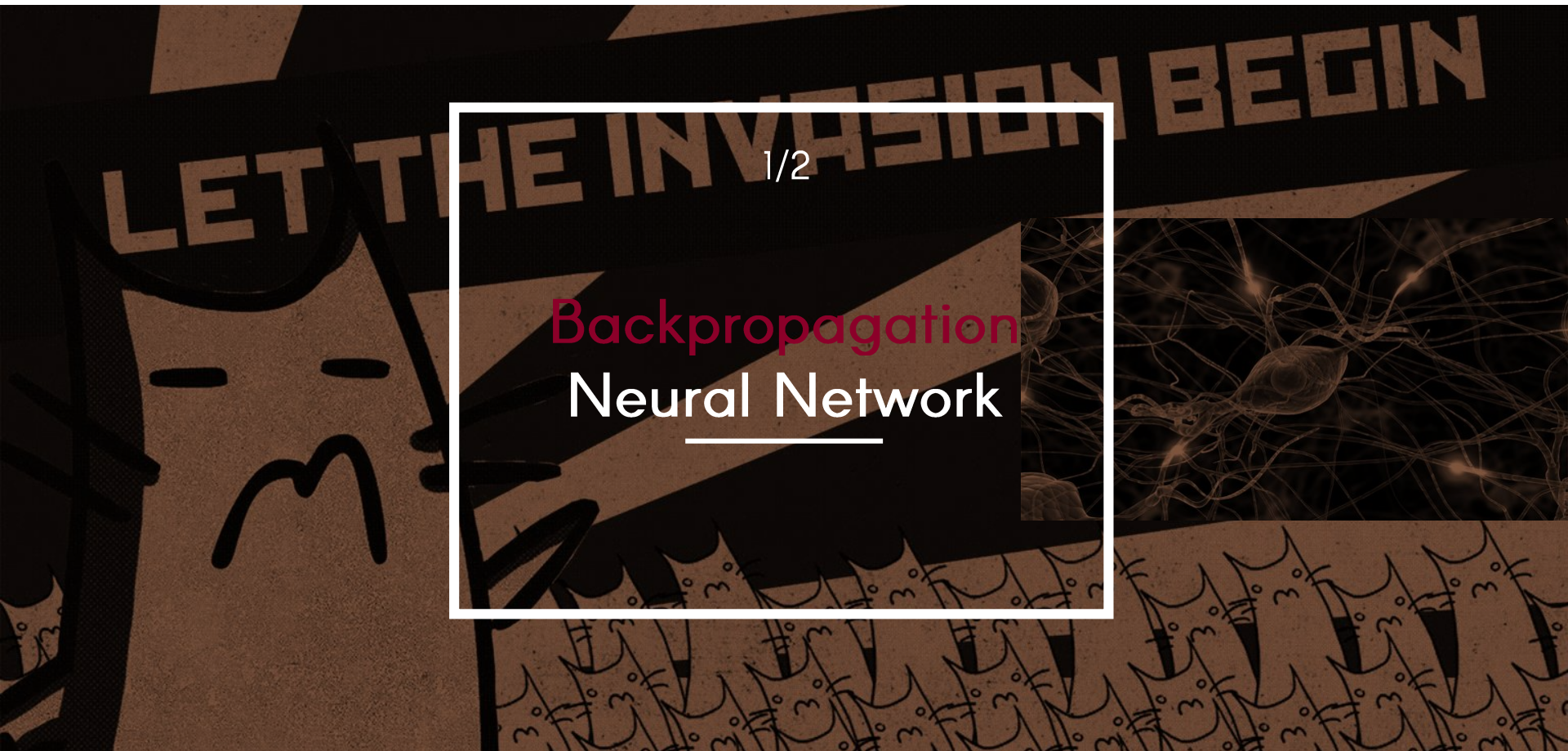


Haedong Kim



1/2

Backpropagation Neural Network



Def.	<ul style="list-style-type: none">• Backpropagation: a way of computing gradients of functions (∇f)
Problem statement	<ul style="list-style-type: none">▪ By applying chain rule recursively▪ In a computational way
Purposes	<ol style="list-style-type: none">1. Perform a parameter update2. Other useful purposed<ul style="list-style-type: none">• Visualization• Interpreting neural networks

- We deal with compounded functions in neural networks

- Multiple hidden layers

1. Support Vector Machine (SVM)

(From the lecture note of Professor Seok, H., Introduction to Machine Learning, Korea Univ., 2016)

- Convex optimization (single global optimum)

2. Neural Network (NN)

- Fixing the number of basis functions in advance
- Allow them to be adaptive in *parametric form*
- Non-convex optimization -> local minimums

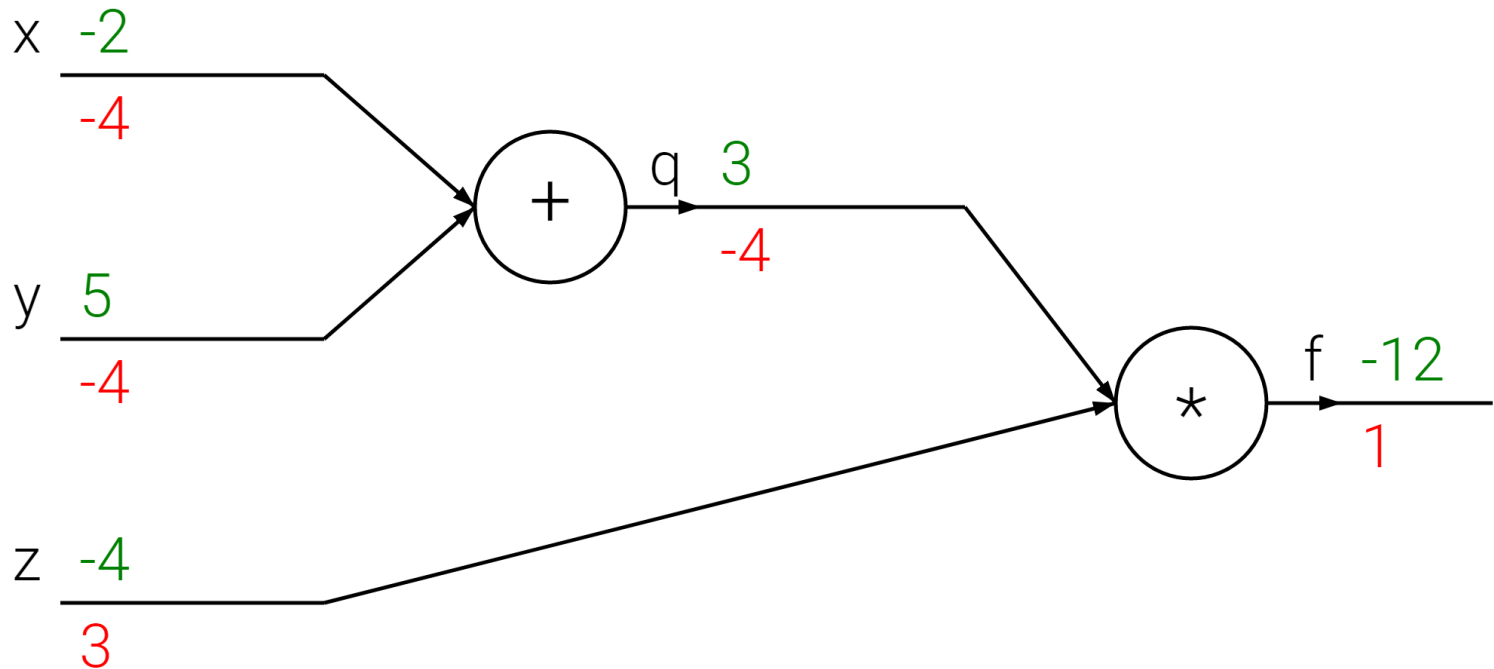
Computational graph

source: <http://cs231n.stanford.edu/syllabus.html>

e.g.

- $f(x, y, z) = (x + y)z$

Computational graph of the f



Goal

- $\frac{\partial f}{\partial z}, \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}$

Computational graph

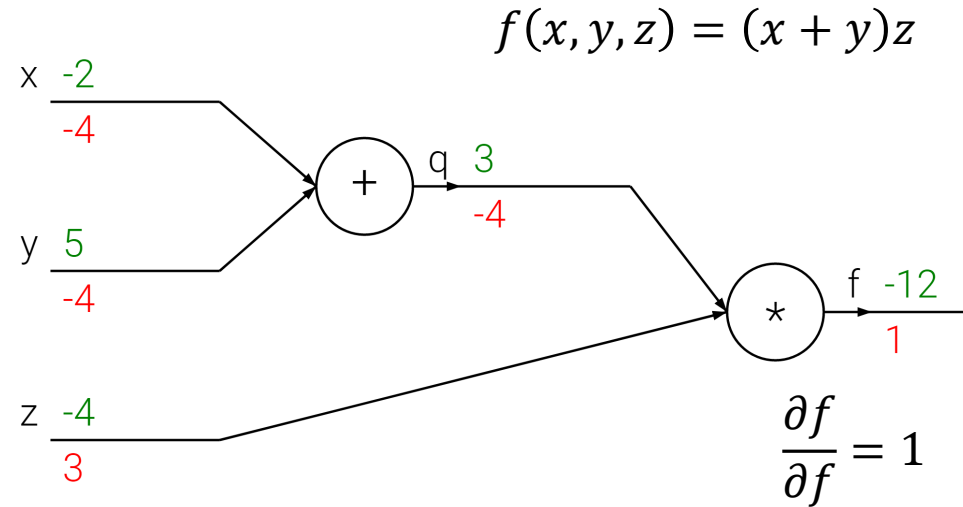
source: <http://cs231n.stanford.edu/syllabus.html>

Multiplication

- $f(x, y) = xy$
 - $\frac{\partial f}{\partial x} = y, \frac{\partial f}{\partial y} = x$

Summation

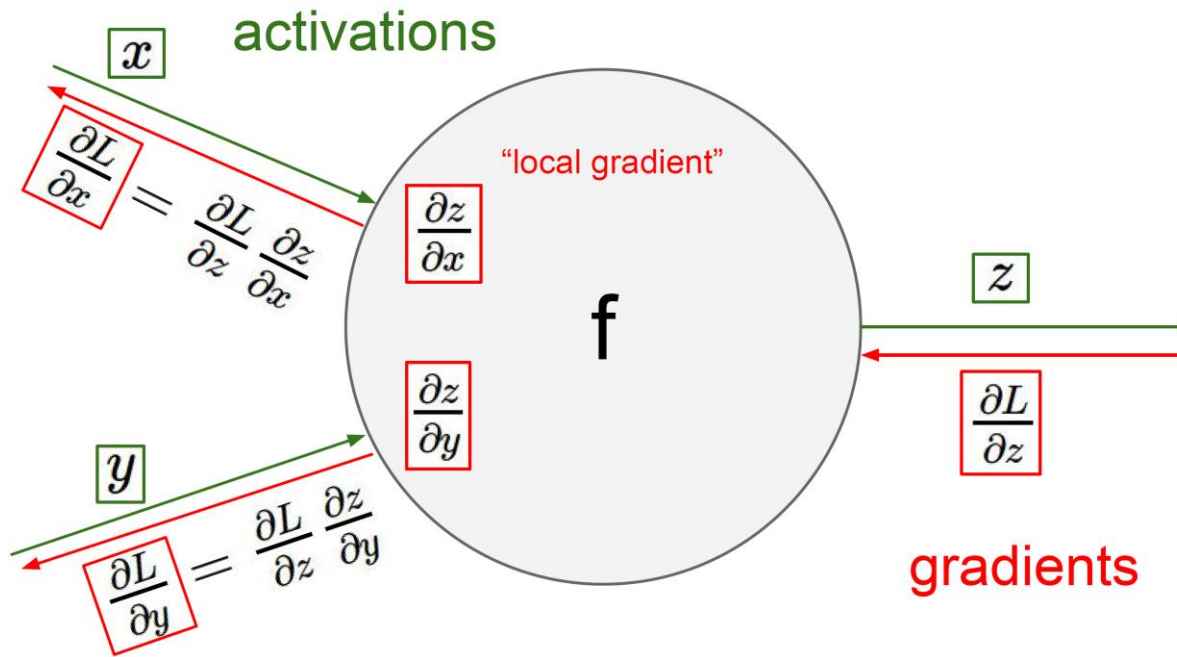
- $f(x, y) = x + y$
 - $\frac{\partial f}{\partial x} = 1, \frac{\partial f}{\partial y} = 1$
- Let $x + y = q$ and $f(x, y, z) = qz$
- $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}, \frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$



Local property of backpropagation

source: <http://cs231n.stanford.edu/syllabus.html>

- Backpropagation is a local process!



- Backpropagation can be thought of as gates communicating to each other

Local property of backpropagation

source: <http://cs231n.stanford.edu/syllabus.html>

e. g.
Sigmoid
function

- $f(\mathbf{w}, \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}} = \frac{1}{1+e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$

Local
gradients of
the sigmoid

- $f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -\frac{1}{x^2}$

- $f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$

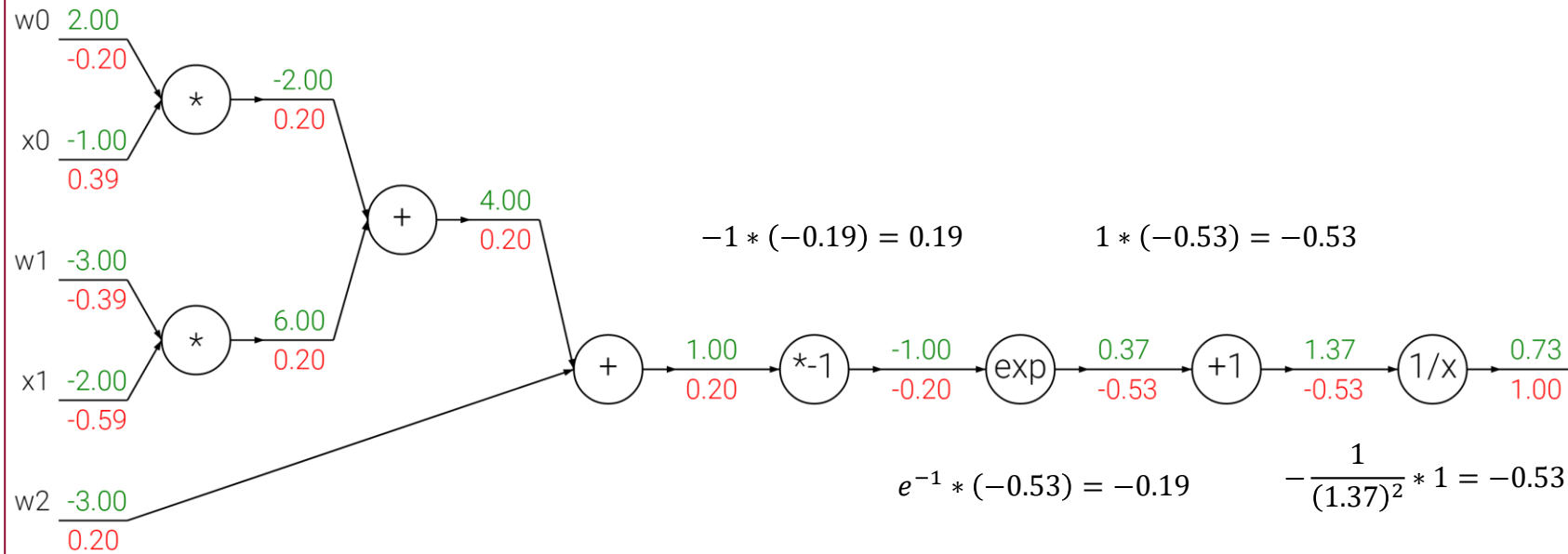
- $f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$

- $f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$

Local property of backpropagation

source: <http://cs231n.stanford.edu/syllabus.html>

Computational graph of the sigmoid



$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

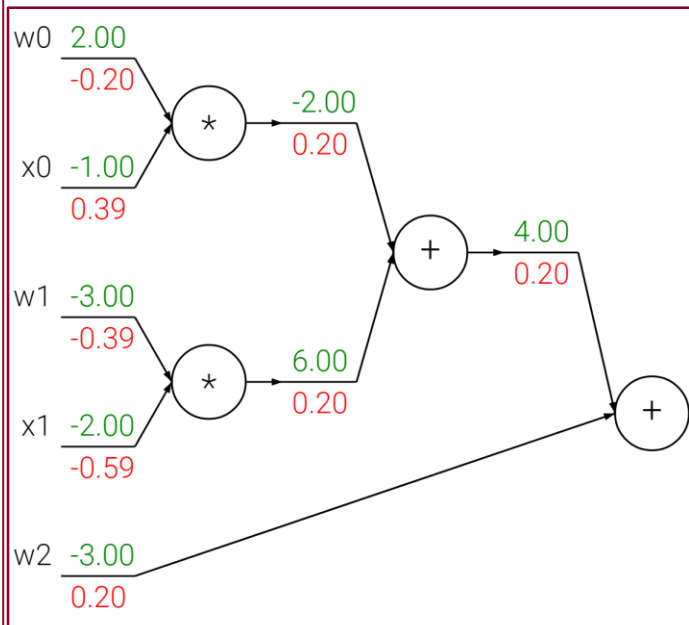
$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

Modularity

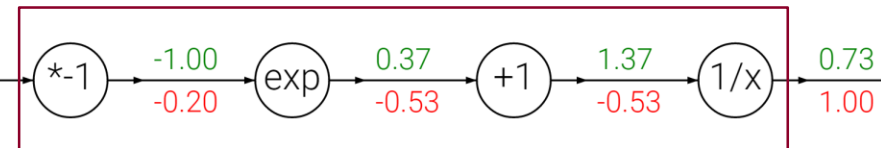
source: <http://cs231n.stanford.edu/syllabus.html>

Merge the graph



← Input of the sigmoid

sigmoid



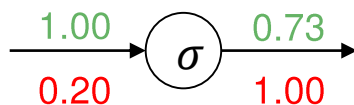
Sigmoid fn.
 $\sigma(x)$

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

Derivative of
the sigmoid

- $\frac{d\sigma(x)}{dx} = (1 - \sigma(x))\sigma(x)$

Sigmoid
gate



$$(1 - 0.73) * 0.73 * 1 = 0.2$$

Patterns in backward flow

source: <http://cs231n.stanford.edu/syllabus.html>

Multiplication gate

- Gradient swapper

- $f(x, y) = xy$

- $\frac{\partial f}{\partial x} = y, \frac{\partial f}{\partial y} = x$

Add gate

- Gradient distributor

- $f(x, y) = x + y$

- $\frac{\partial f}{\partial x} = 1, \frac{\partial f}{\partial y} = 1$

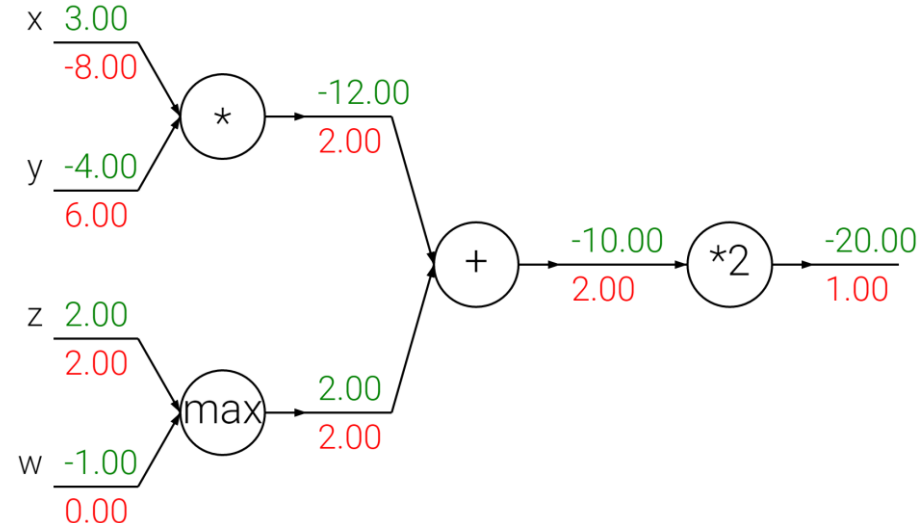
Max gate

- Gradient router

- $f(x, y) = \max(x, y)$

- $\frac{\partial f}{\partial x} = 1 \ (x \geq y), 0 \ o/w$

- $\frac{\partial f}{\partial y} = 1 \ (y \geq x), 0 \ o/w$



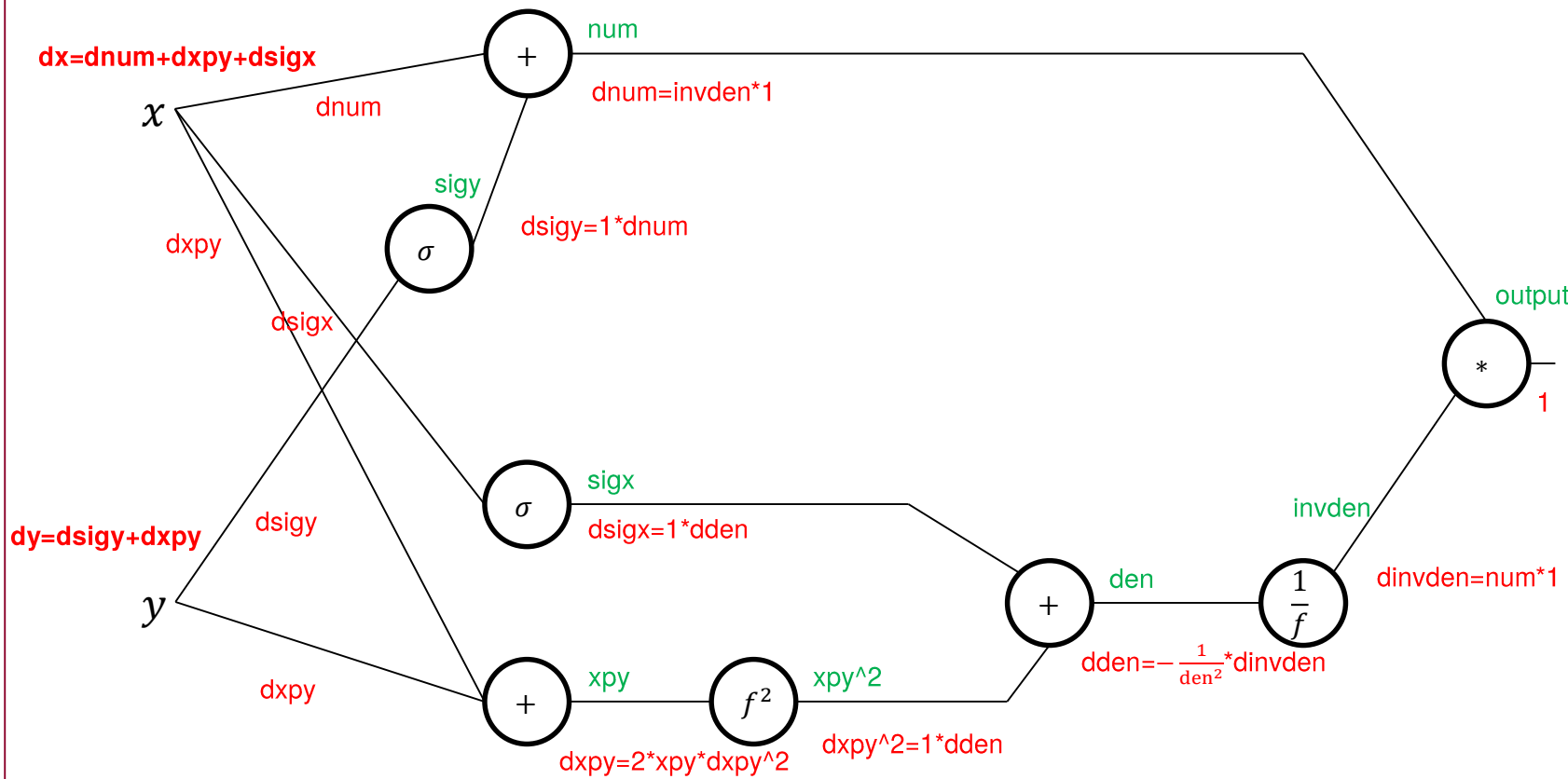
Another example

source: <http://cs231n.stanford.edu/syllabus.html>

Any function

- $$f(x, y) = \frac{x + \sigma(y)}{\sigma(x) + (x + y)^2}$$
 (It's useless in practice)

Computational graph



dx, dy

- add up gradients at forks

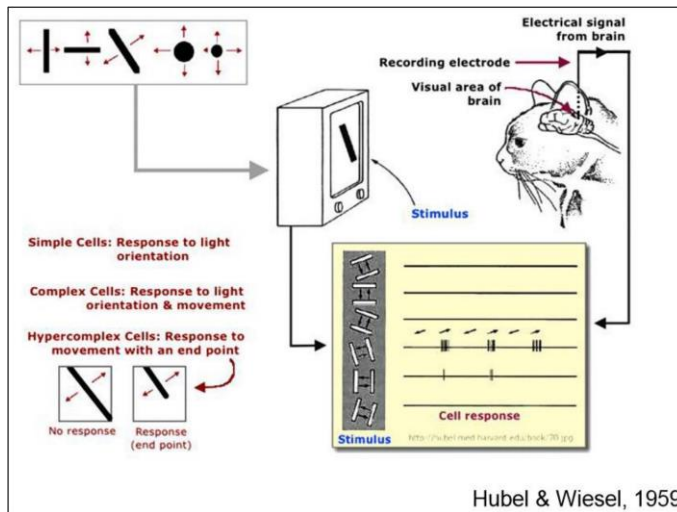
Problem	<ul style="list-style-type: none">• Neural nets will be very large<ul style="list-style-type: none">▪ Its gradient also will be very complex
Modularity	<ul style="list-style-type: none">• Recursive application of the chain rule will save us!<ul style="list-style-type: none">▪ Local property
Computational way	<ul style="list-style-type: none">• It turns out that we don't need explicit expressions of gradients<ul style="list-style-type: none">▪ Intermediate values▪ Staged computation

1/2

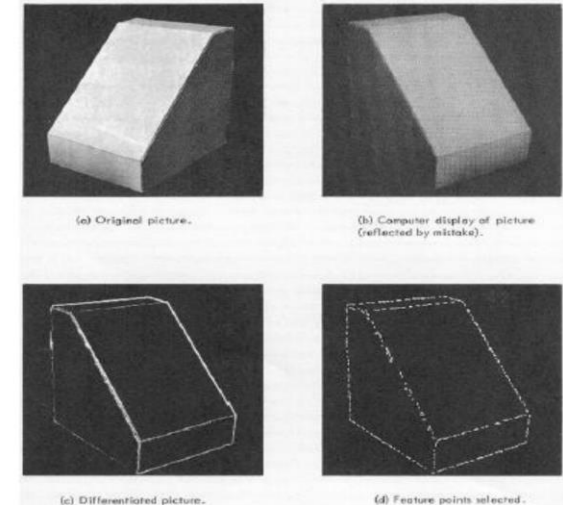
Backpropagation Neural Network

- “The area of Neural Networks has originally been primarily inspired by the goal of modeling biological neural systems, but has since diverged and become a matter of engineering...” (*Andrej Karpathy*)

Revisit



Block
world
Larry Roberts,
1963



Parallel
processing

- Parallelism** is the main difference of the human brain from a computer

SIMD

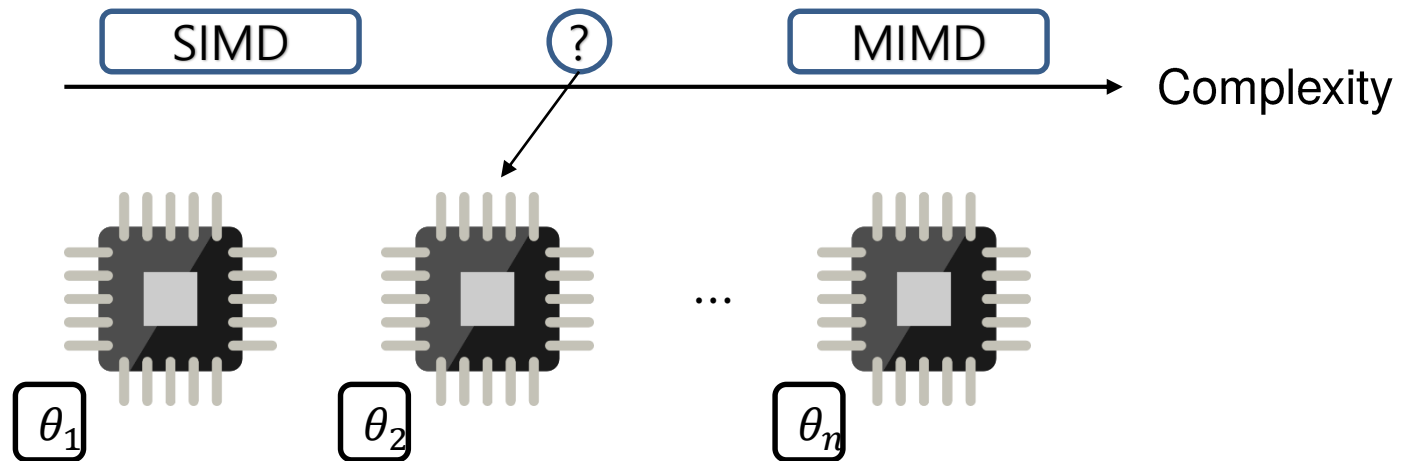
- Two paradigms for *parallel processing*
- Single Instruction Multiple Data (SIMD)
 - All processors execute the same instruction but on different pieces of data

MIMD

- Multiple Instruction Multiple Data (MIMD)
 - Different processors may execute different instructions on different data

Parallel processing

Alpaydin, E. (2014). Introduction to machine learning: MIT press.



- Each processor is a fixed function and execute same instruction as SIMD
- By loading different parameters into a local memory, they can be doing different thing
- The whole operation can be **distributed** over such processors

Distributed representation



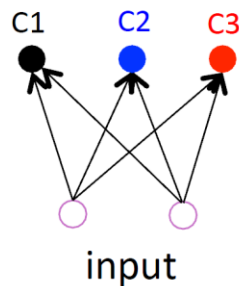
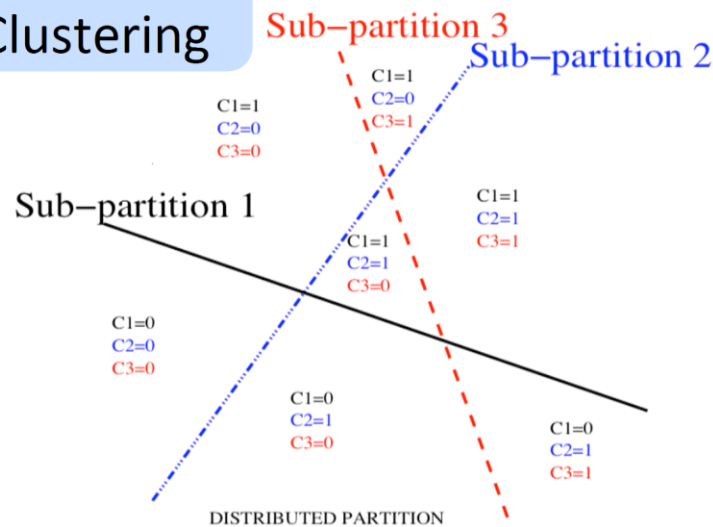
Bengio, B. (2015). Deep learning: Theoretical Motivations. Retrieved from http://videolectures.net/deeplearning2015_bengio_theoretical_motivations/

Curse of dim.

- Distributed representation allows us bypass the curse of dimensionality

The need for dist. rep.

Multi-Clustering



Non-mutually exclusive features/ attributes create a combinatorially large set of distinguishable configurations

of distinguishable regions is linear in # of parameters

of distinguishable regions grows almost exponentially with # of parameters

Generalize non-locally to never seen regions

Architecture

- Depth and width

Universal approximation thm.

- A feedforward network with a single layer can approximate any Borel measurable function from one finite-dim. space to another, provided that the network is given enough hidden units

(Hornik, K., Stinchcombe, M., & White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5), 359-366.)

Problems

1. Doesn't know appropriate parameters (Number of hidden units)
2. Overfitting

Practical solution

- Sufficient capability & Regularization



Q&A