

MAE 688: Machine Learning for Mechanical Engineers

Lecture 2- Introduction to Deep Learning



Dr. Bing Dong

Director, Built Environment Science and Technology (BEST) Lab

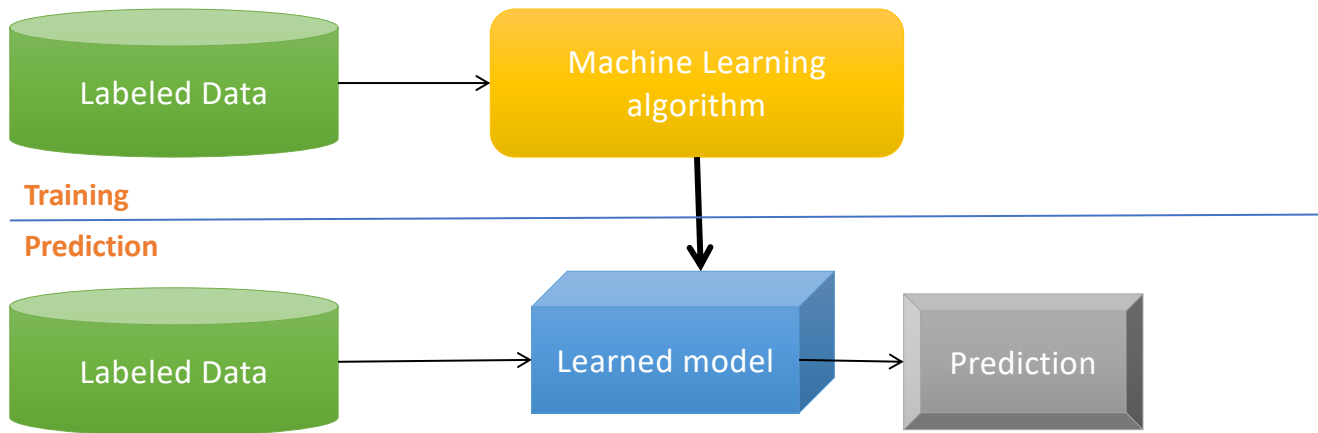
Traugott Professor

Mechanical and Aerospace Engineering

Syracuse University

What is Machine Learning?

Machine learning is a field of computer science that gives computers the ability to **learn without being explicitly programmed**



Methods that can learn from and make predictions on data

Types of Machine Learning

Supervised: Learning with a **labeled training** set

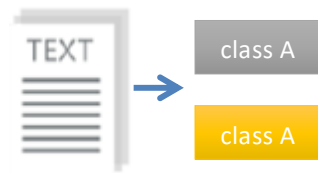
Example: email *classification* with already labeled emails

Unsupervised: Discover **patterns** in **unlabeled** data

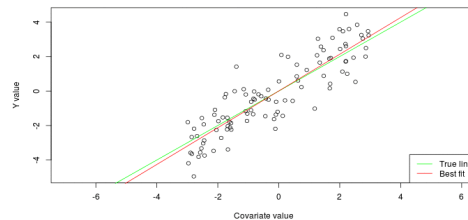
Example: *cluster* similar documents based on text

Reinforcement learning: learn to **act** based on **feedback/reward**

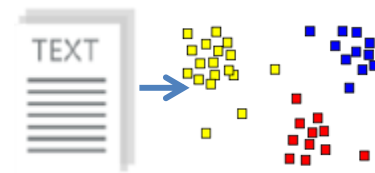
Example: learn to play Go, reward: *win or lose*



Classification



Regression



Clustering

Anomaly Detection
Sequence labeling

...

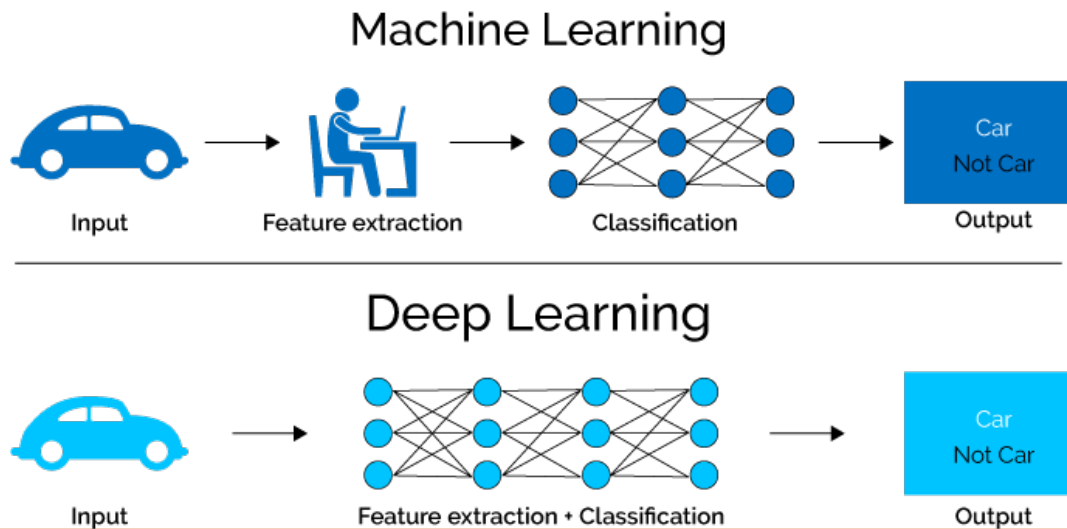
<http://mbjoseph.github.io/2013/11/27/measure.html>

What is Deep Learning?

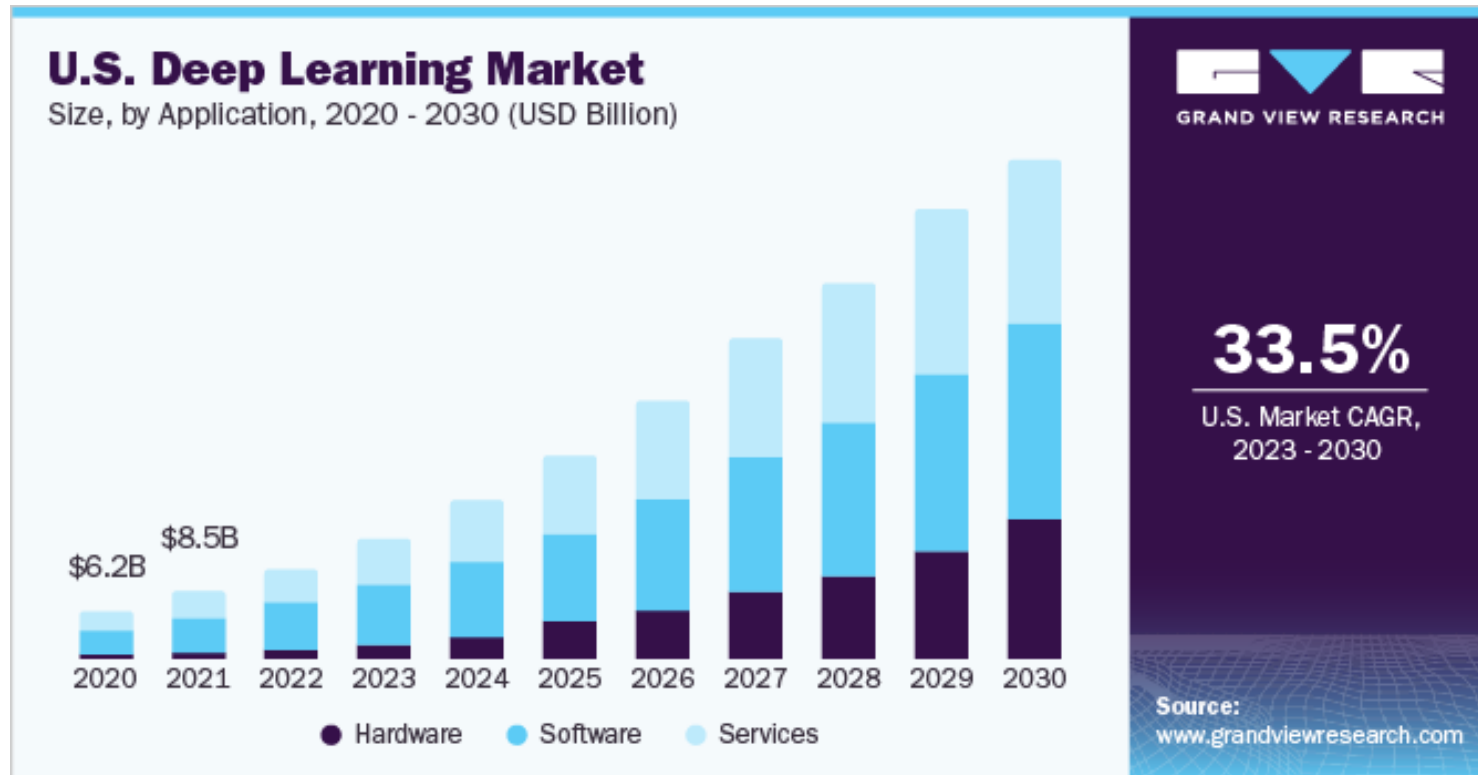
A machine learning subfield of learning **representations** of data. Exceptional effective at **learning patterns**.

Deep learning algorithms attempt to learn (multiple levels of) representation by using a **hierarchy of multiple layers**

If you provide the system **tons of information**, it begins to understand it and respond in useful ways.



Deep learning attracts lots of attention



The History of Deep Learning

1958: Perceptron (linear model)

1969: Perceptron has limitation

1980s: Multi-layer perceptron

Do not have significant difference from DNN today

1986: Backpropagation

Usually more than 3 hidden layers is not helpful

1989: 1 hidden layer is “good enough”, why deep?

2006: RBM initialization

2009: GPU

2011: Start to be popular in speech recognition

2012: win ILSVRC image competition

2015.2: Image recognition surpassing human-level performance

2016.3: Alpha GO beats Lee Sedol

2016.10: Speech recognition system as good as humans

Hinton, G.E., 1986, August. Learning distributed representations of concepts. In *Proceedings of the eighth annual conference of the cognitive science society* (Vol. 1, p. 12).

Hinton, G.E. and Salakhutdinov, R.R., 2006. Reducing the dimensionality of data with neural networks. *science*, 313(5786), pp.504-507.

Restricted boltzmann machines

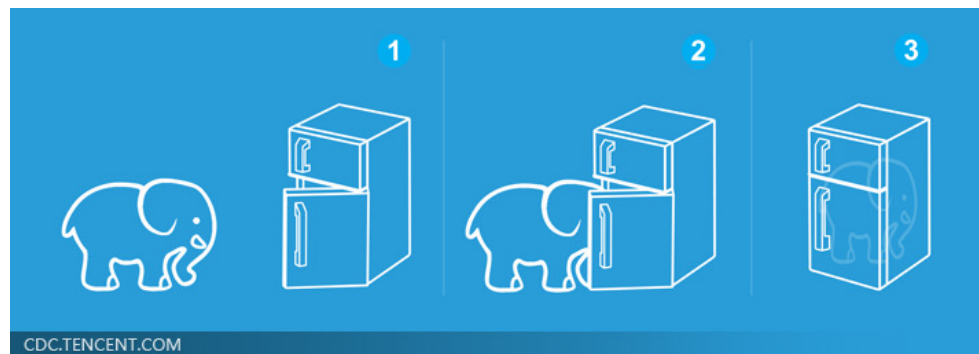
Three "SIMPLE" steps for deep learning



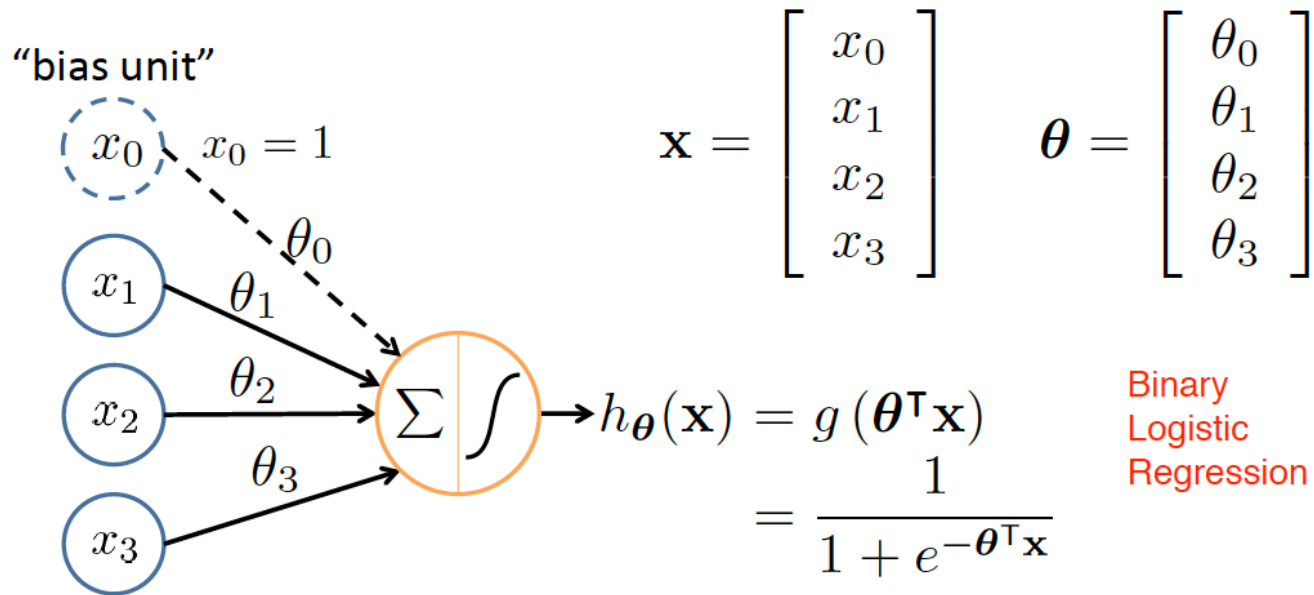
Three "SIMPLE" steps for deep learning



Deep Learning is so simple

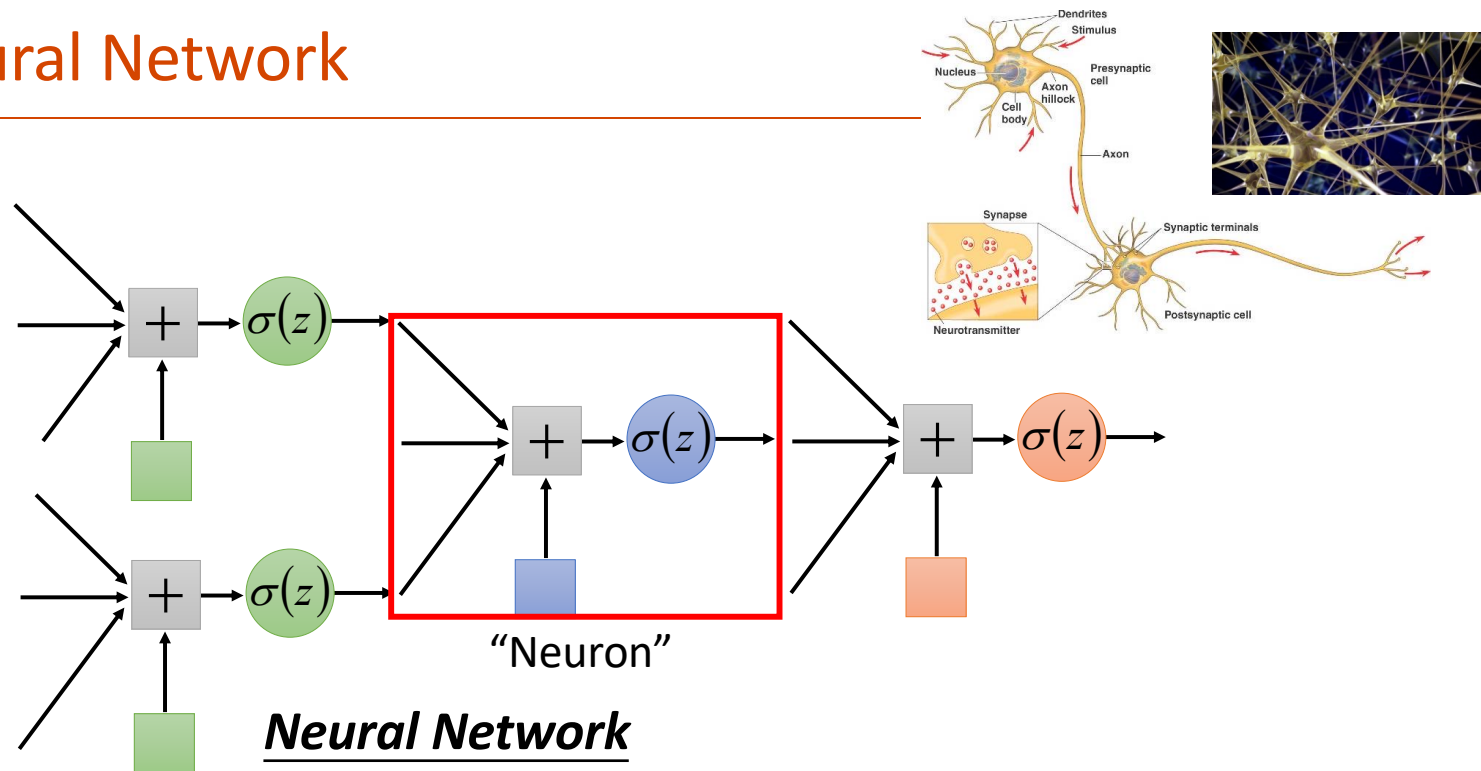


Single Node



Sigmoid (logistic) activation function: $g(z) = \frac{1}{1 + e^{-z}}$

Neural Network

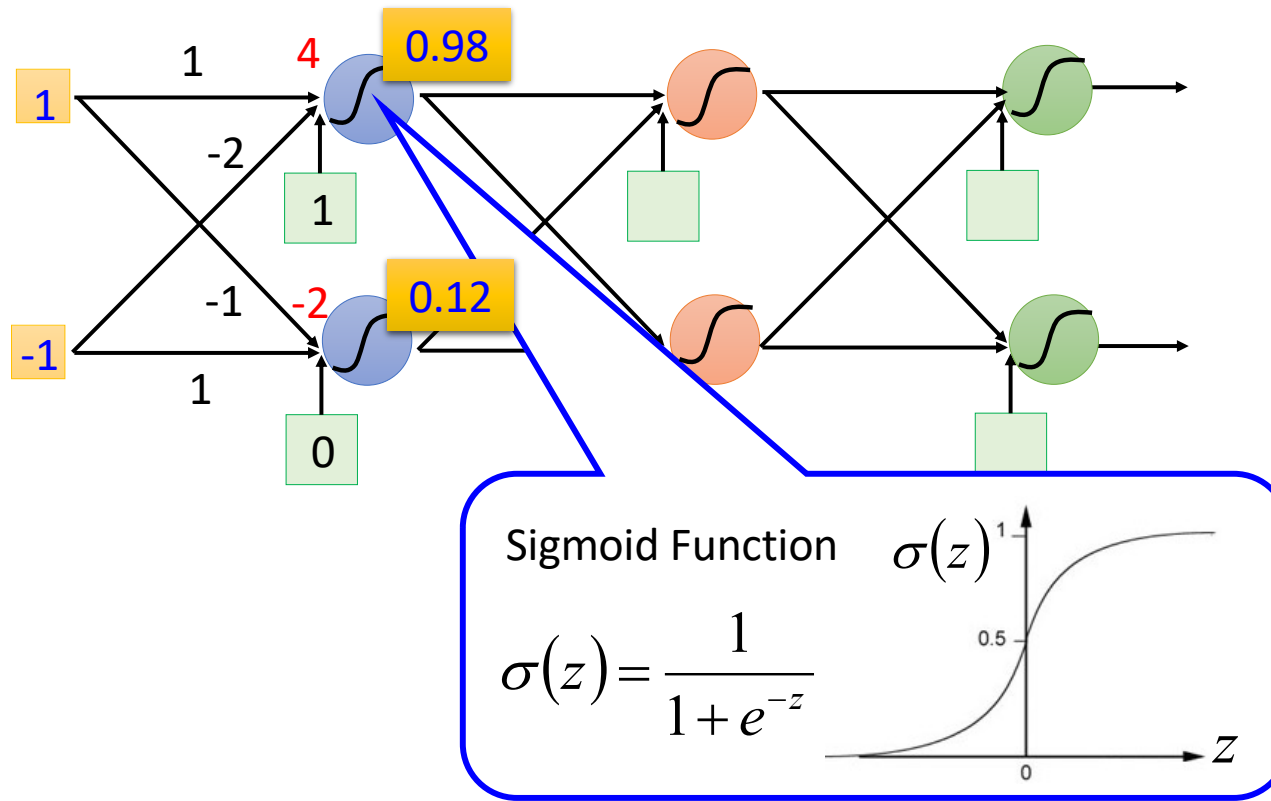


Neural Network

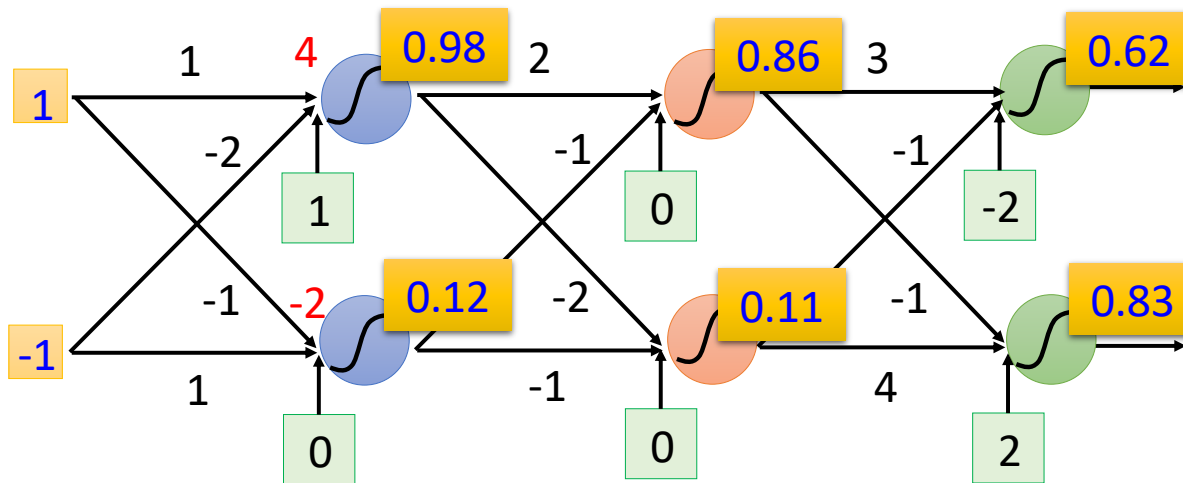
Different connection leads to different network structures

Network parameter θ : all the weights and biases in the "neurons"

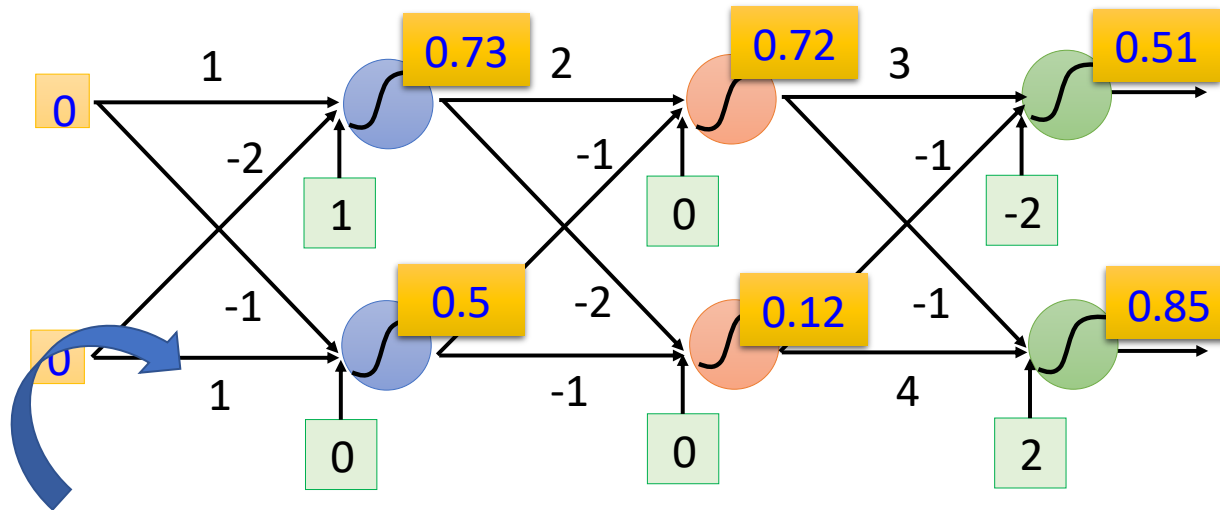
Fully Connect Feedforward Network



Fully Connect Feedforward Network



Fully Connect Feedforward Network

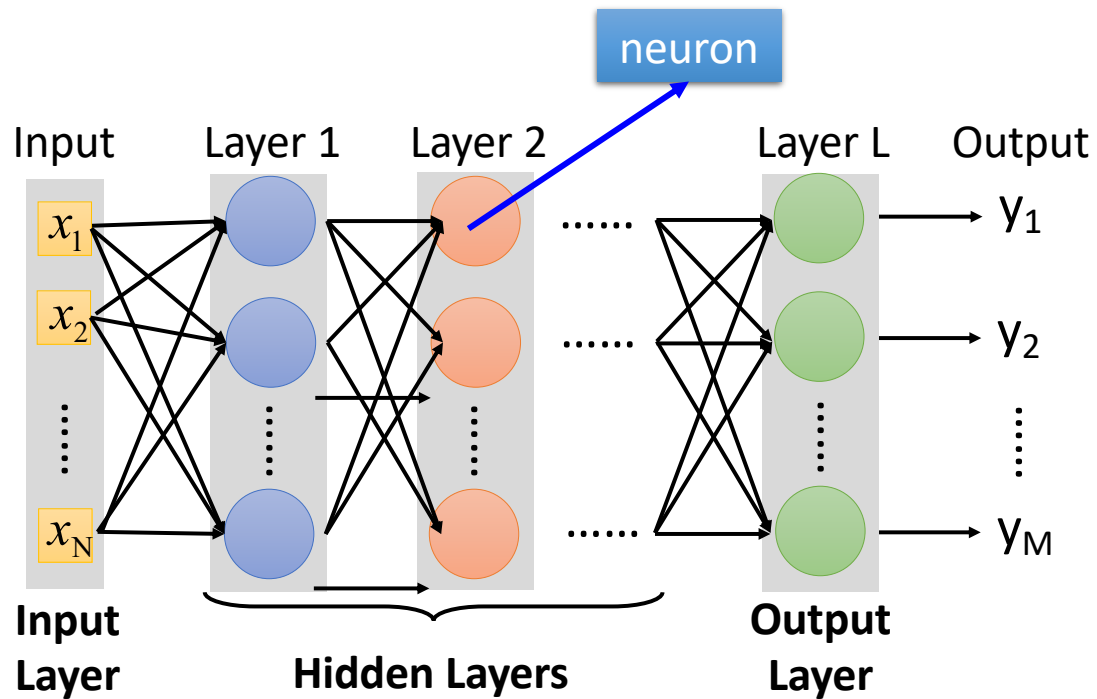


This is a function.
Input vector, output vector

$$f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

Given network structure, define a function set

Fully Connect Feedforward Network

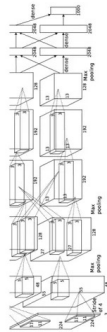


Deep = Many Many hidden la

Large Scale Visual Recognition Challenge (ILSVRC)

16.4%

8 layers



AlexNet (2012)

7.3%

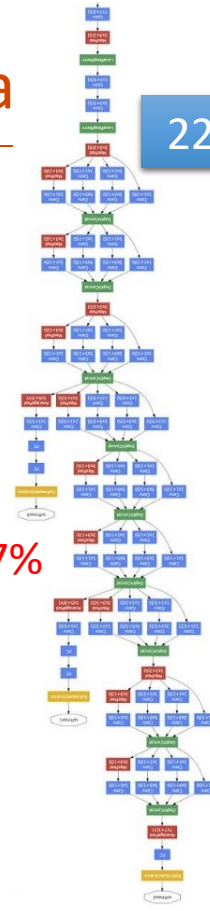
19 layers



VGG (2014)

6.7%

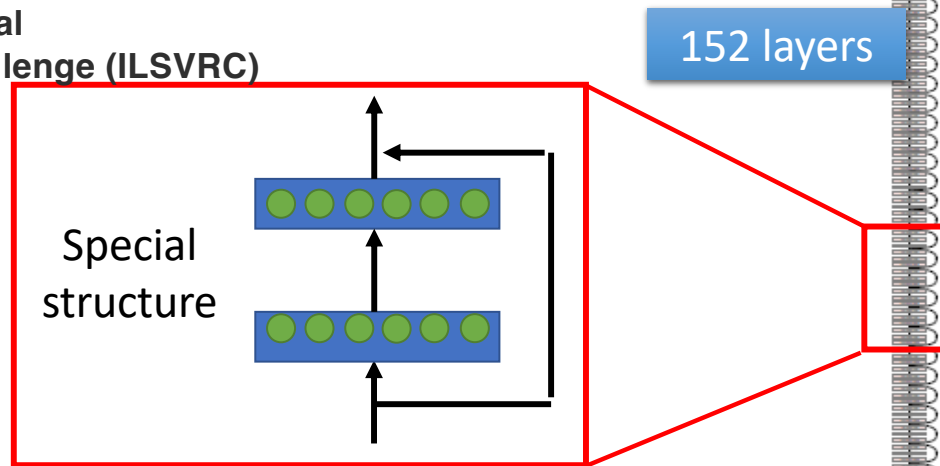
22 layers



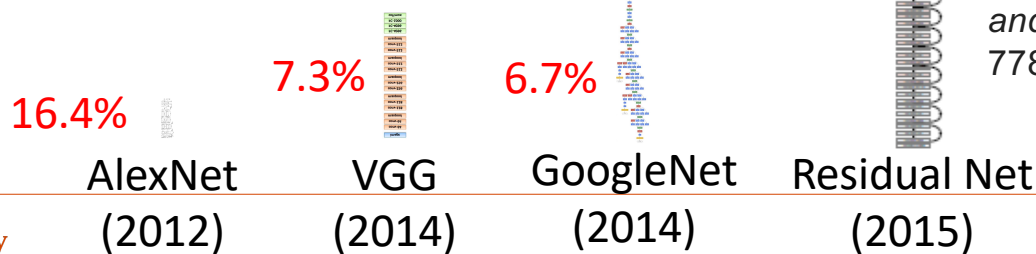
GoogleNet (2014)

Deep = Many Many hidden layers

Large Scale Visual Recognition Challenge (ILSVRC)



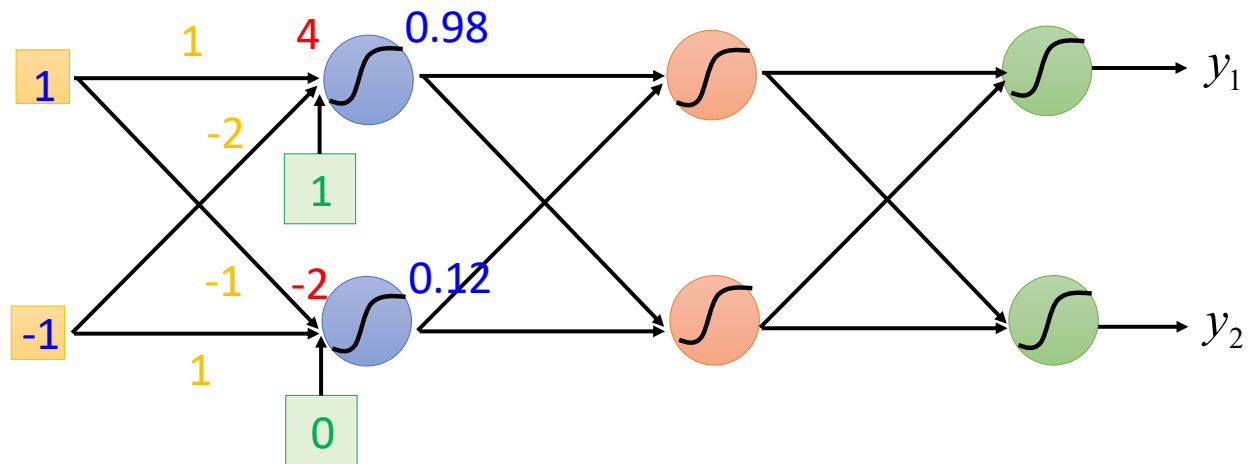
Ref:
<https://www.youtube.com/watch?v=dxB6299gpvl>



He, K., Zhang, X., Ren, S. and Sun, J., 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

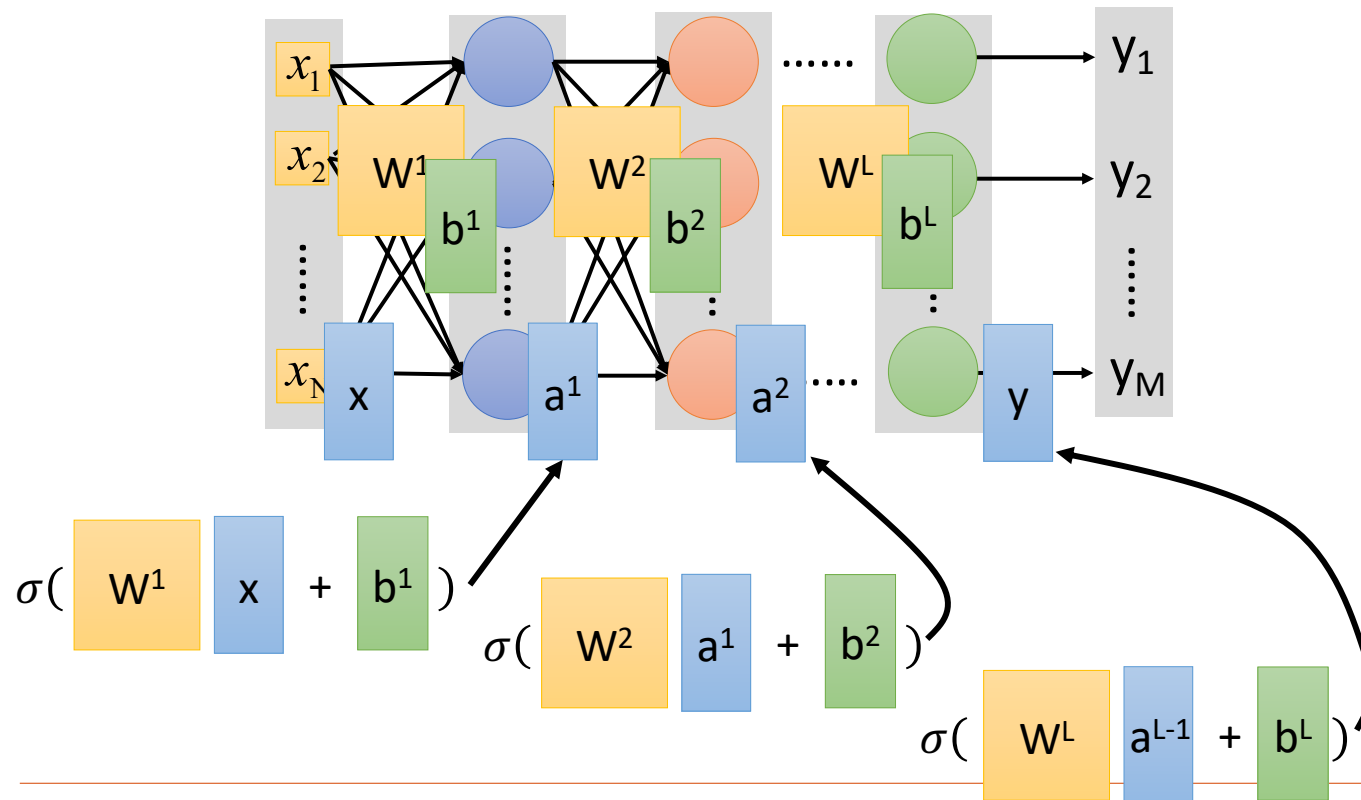
[Cited by 235134](#)

Matrix Operation

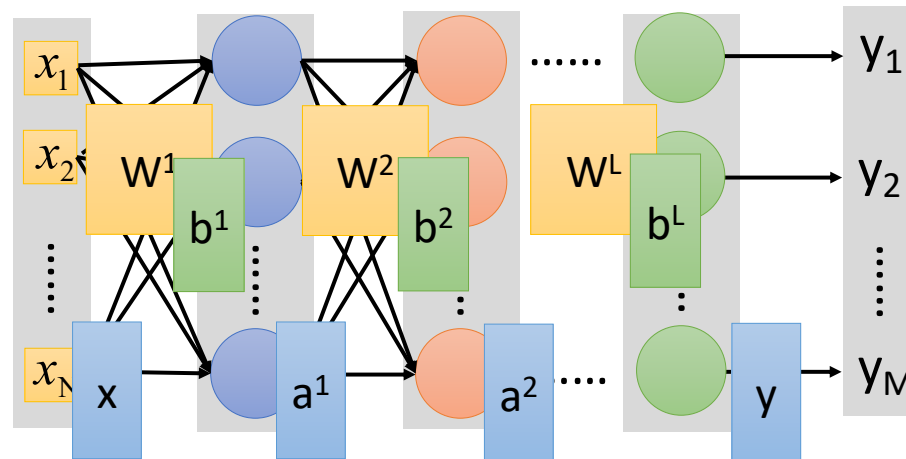


$$\sigma\left(\underbrace{\begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}}_{\begin{bmatrix} 4 \\ -2 \end{bmatrix}} \right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

Neural Network



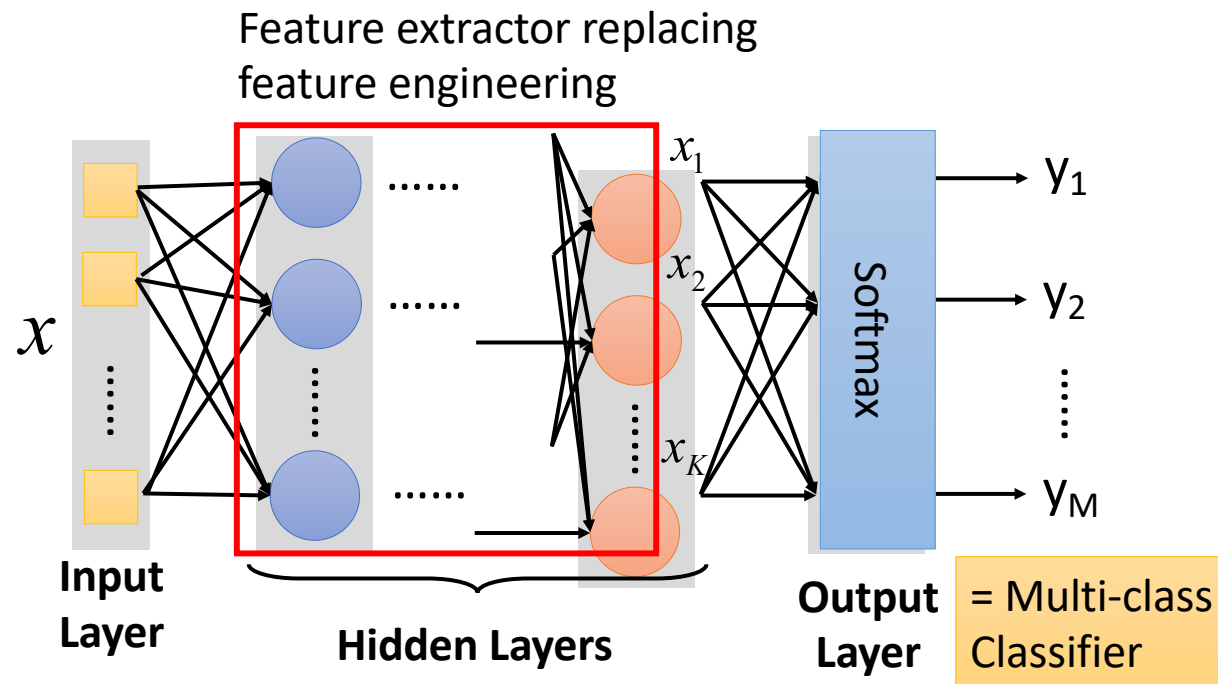
Neural Network



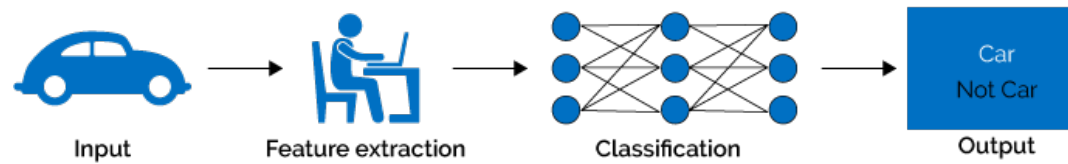
$y = f(x)$ Using parallel computing techniques
 to speed up matrix operation

$$= \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

Output Layer as Multi-Class Classifier



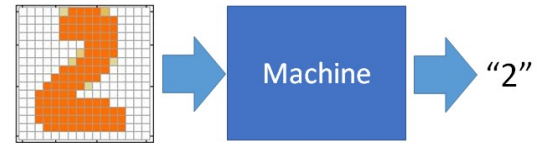
Machine Learning



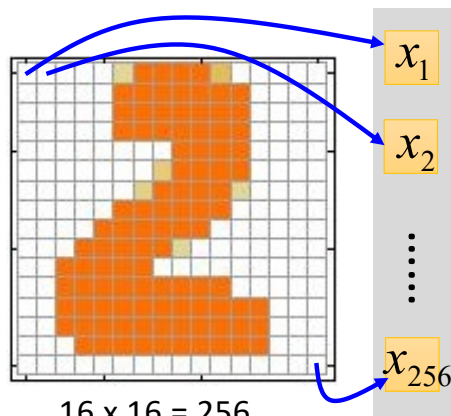
Deep Learning



Example Application



- Input

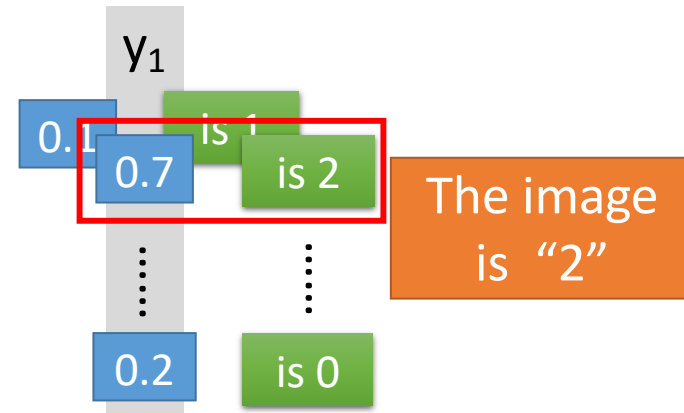


16 x 16 = 256

Ink → 1

No ink → 0

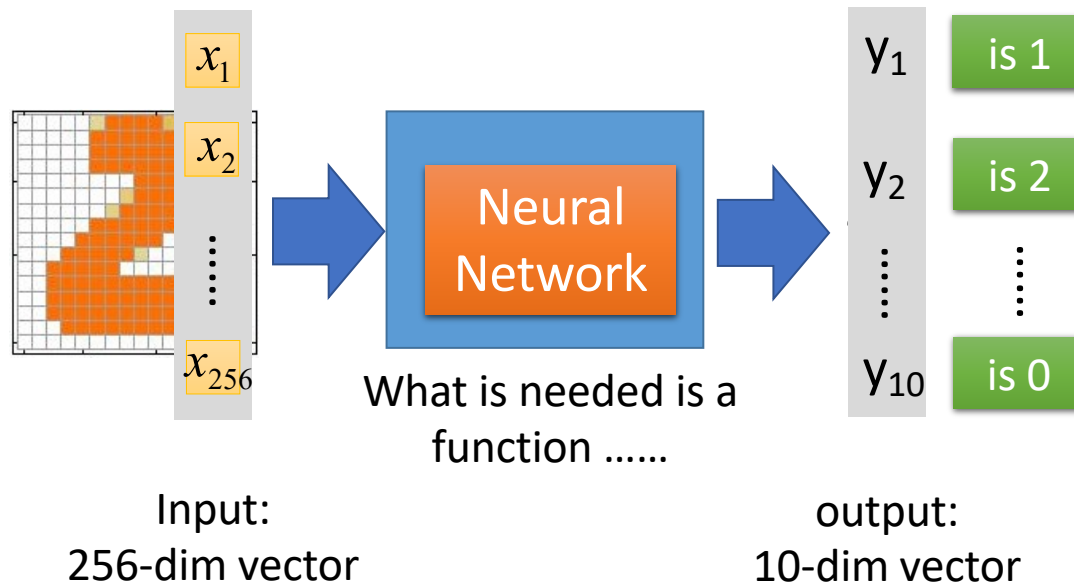
- Output



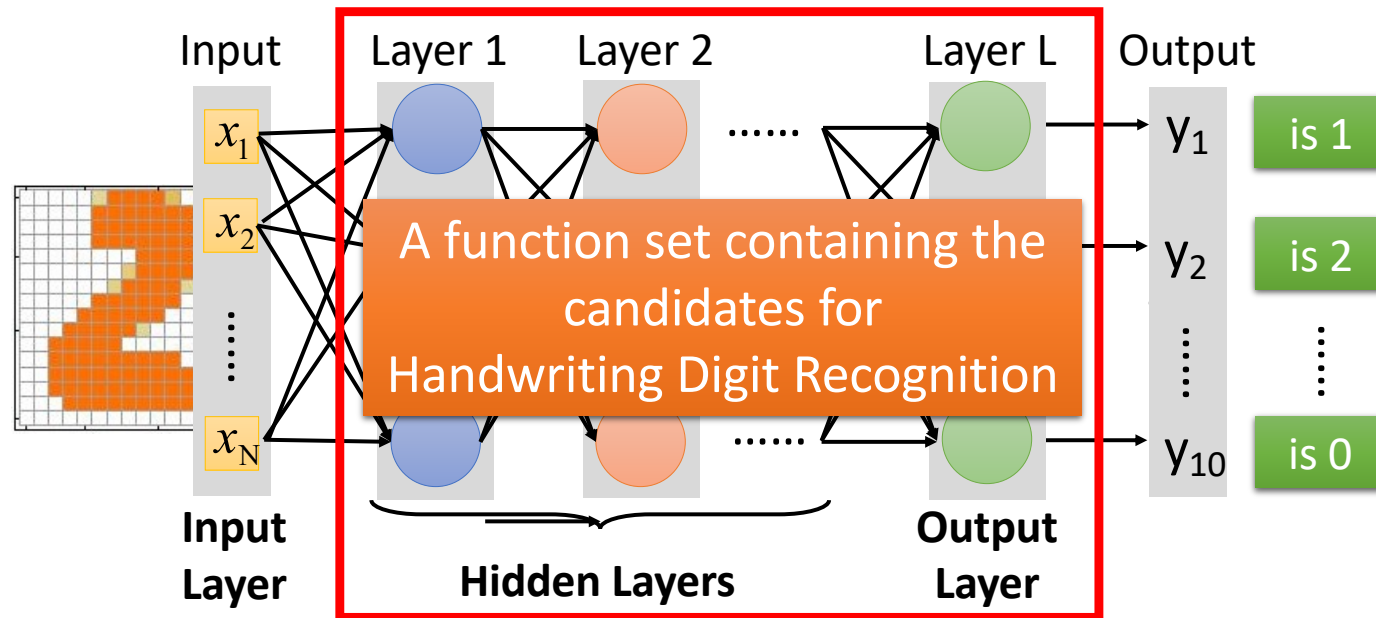
Each dimension represents the confidence of a digit.

Example Application

- Handwriting Digit Recognition

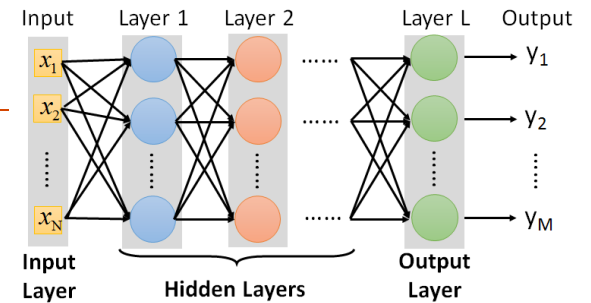


Example Application



You need to decide the network structure to let a good function in your function set.

FAQ



- Q: How many layers? How many neurons for each layer?

Trial and Error

+

Intuition

- Q: Can the structure be automatically determined?
 - E.g. Evolutionary Artificial Neural Networks
- Q: Can we design the network structure?

Convolutional Neural Network (CNN)

Activation Functions

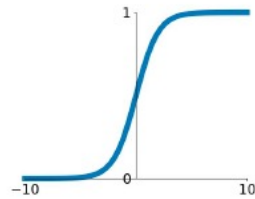
(borrow some slides are from Stanford University)



Activation Functions

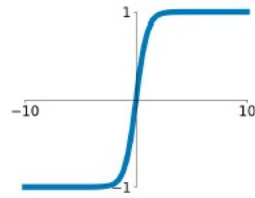
Sigmoid

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



tanh

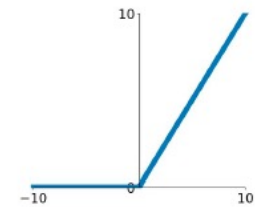
$$\tanh(x)$$



ReLU

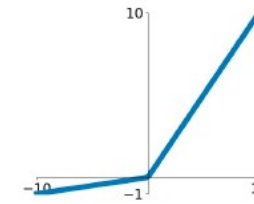
$$\max(0, x)$$

Rectified Linear Unit



Leaky ReLU

$$\max(0.1x, x)$$



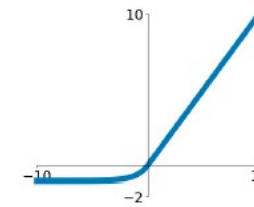
Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

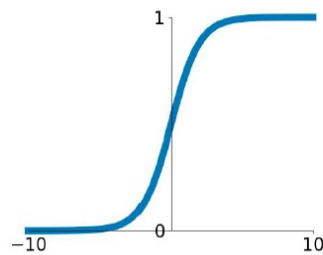
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

Exponential Linear Units



Activation Functions



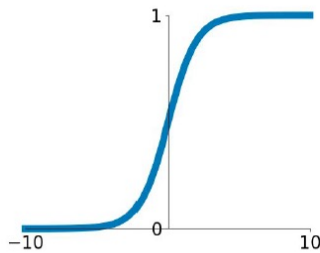
Sigmoid

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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

Activation Functions



Sigmoid

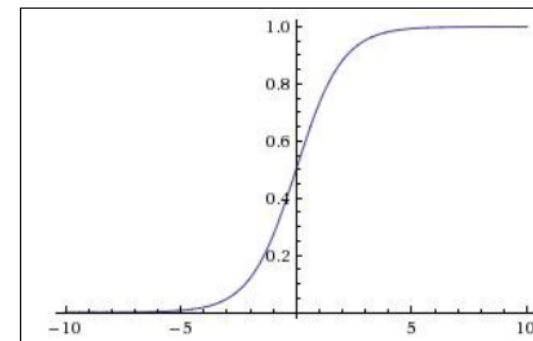
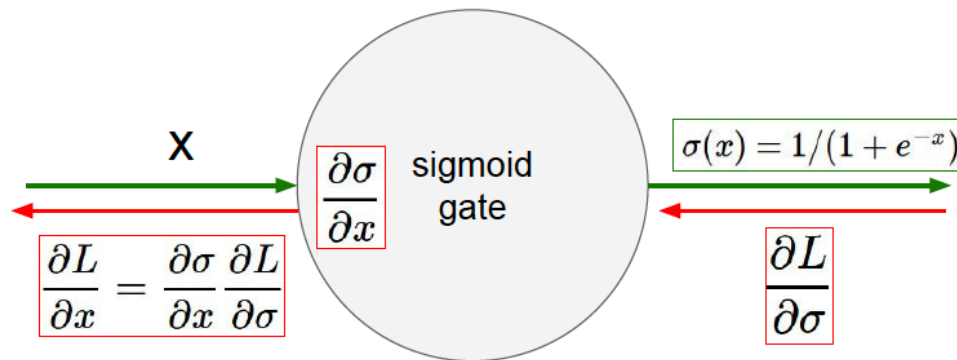
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients

Activation Functions



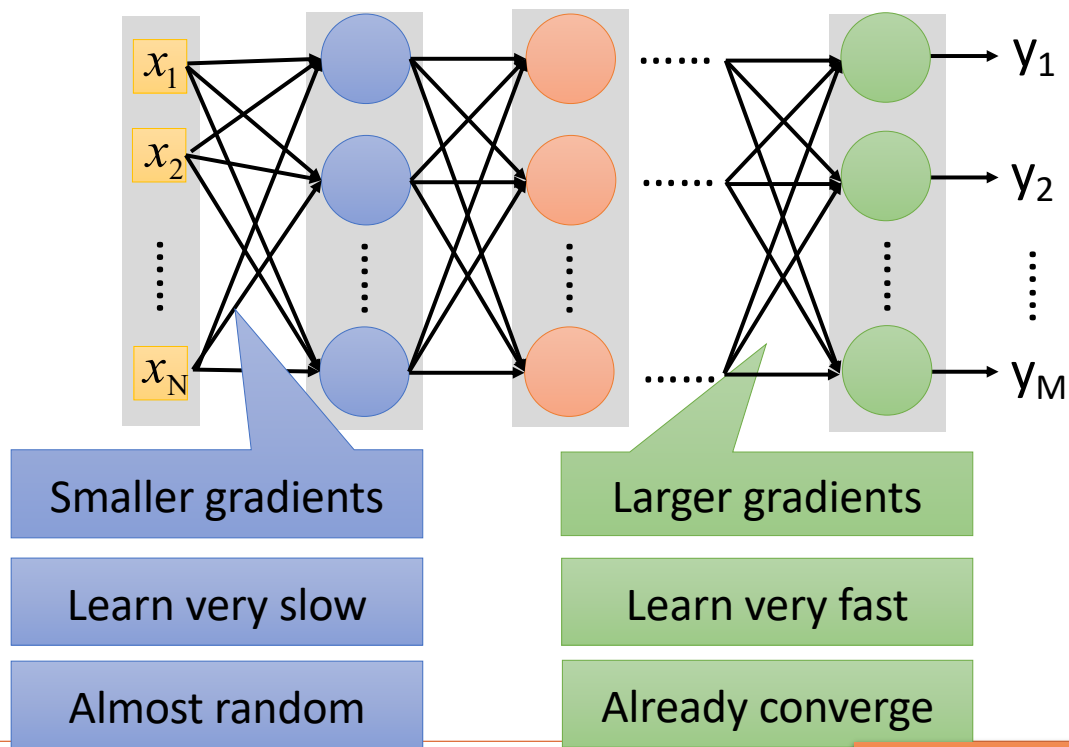
What happens when $x = -10$?

What happens when $x = 0$?

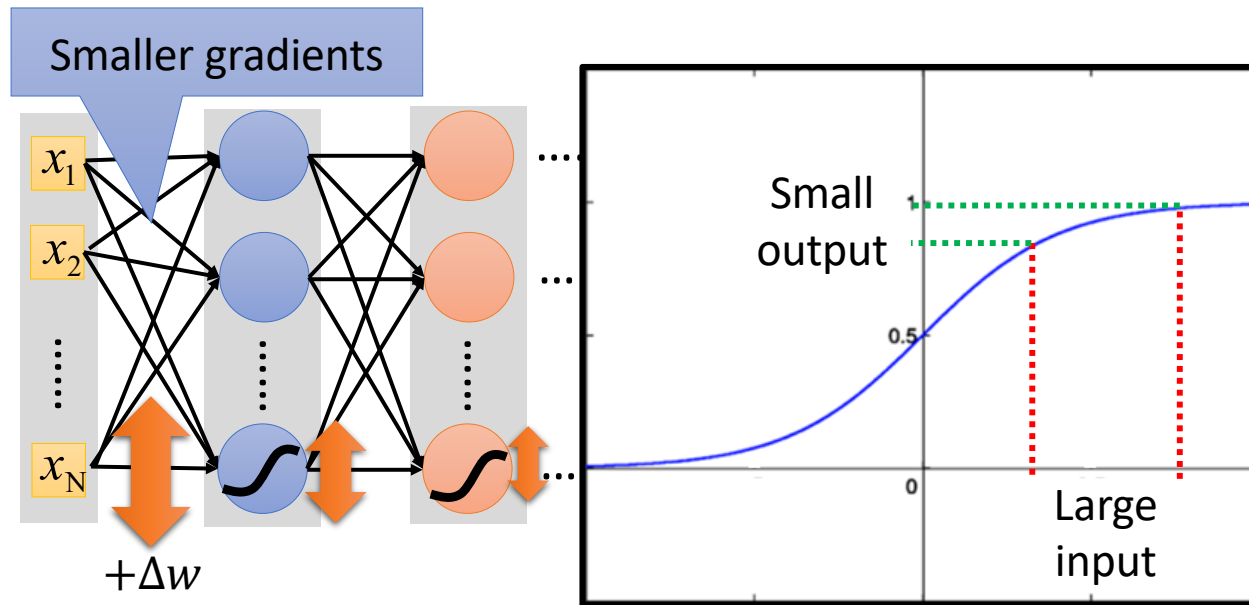
What happens when $x = 10$?

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

Vanishing Gradient Problem



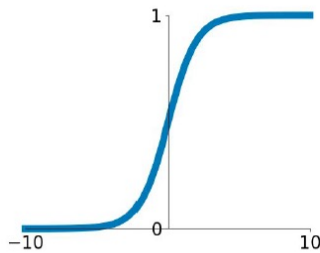
Vanishing Gradient Problem



Intuitive way to compute the derivatives ...

$$\frac{\partial l}{\partial w} \stackrel{?}{=} \frac{\Delta l}{\Delta w}$$

Activation Functions



Sigmoid

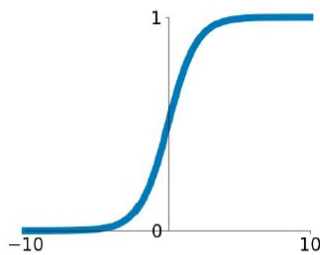
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered

Activation Functions



Sigmoid

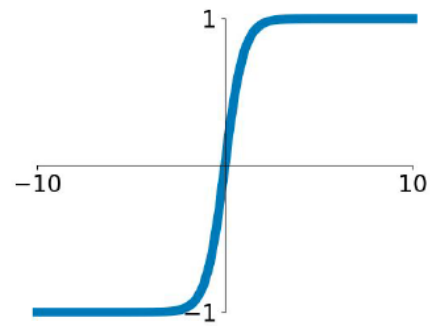
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit compute expensive

Activation Functions

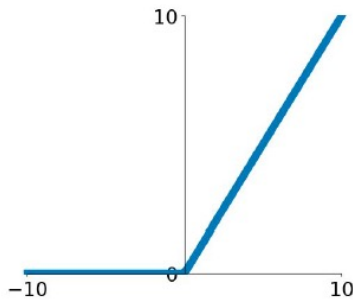


tanh(x)

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(

[LeCun et al., 1991]

Activation Functions



ReLU
(Rectified Linear Unit)

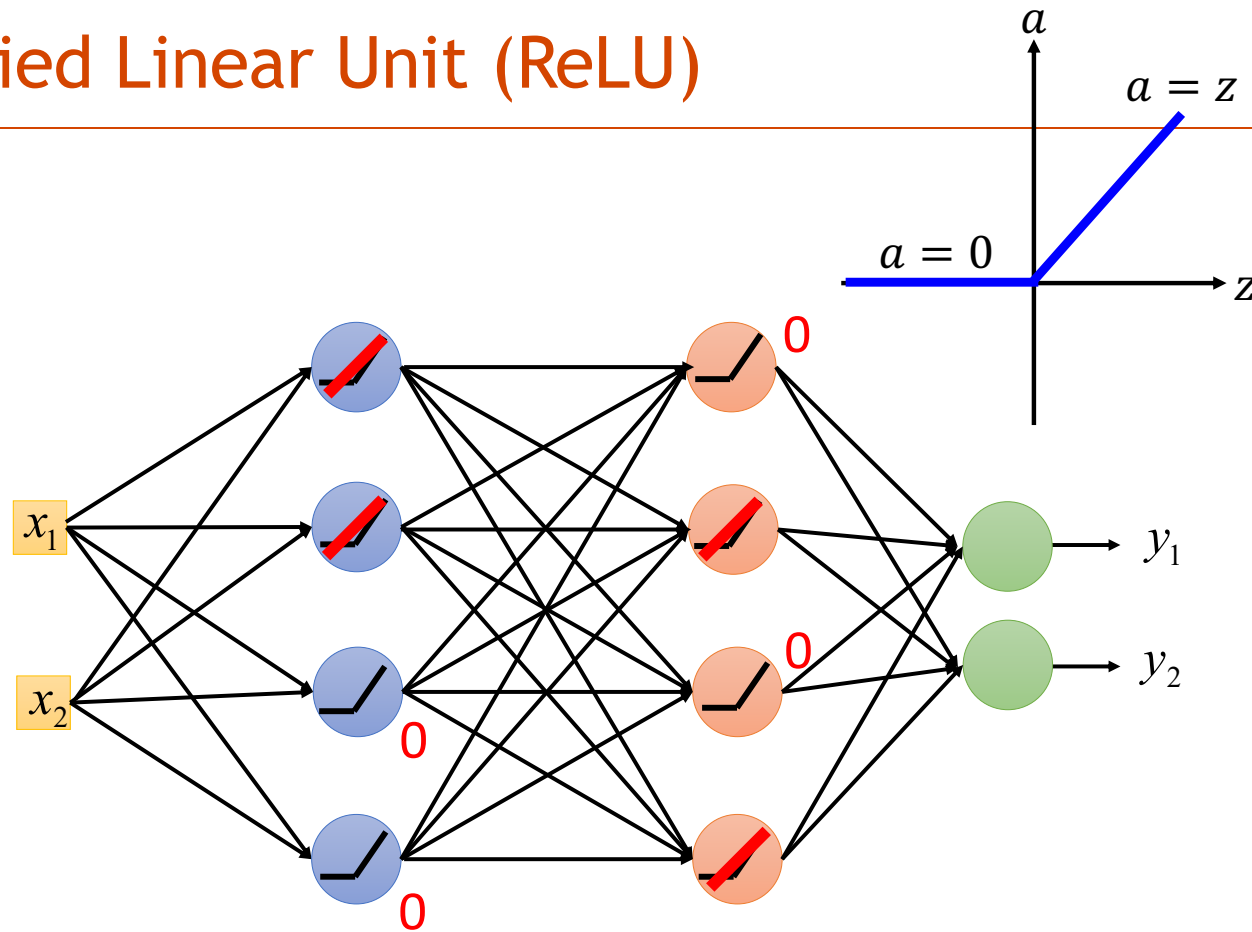
- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

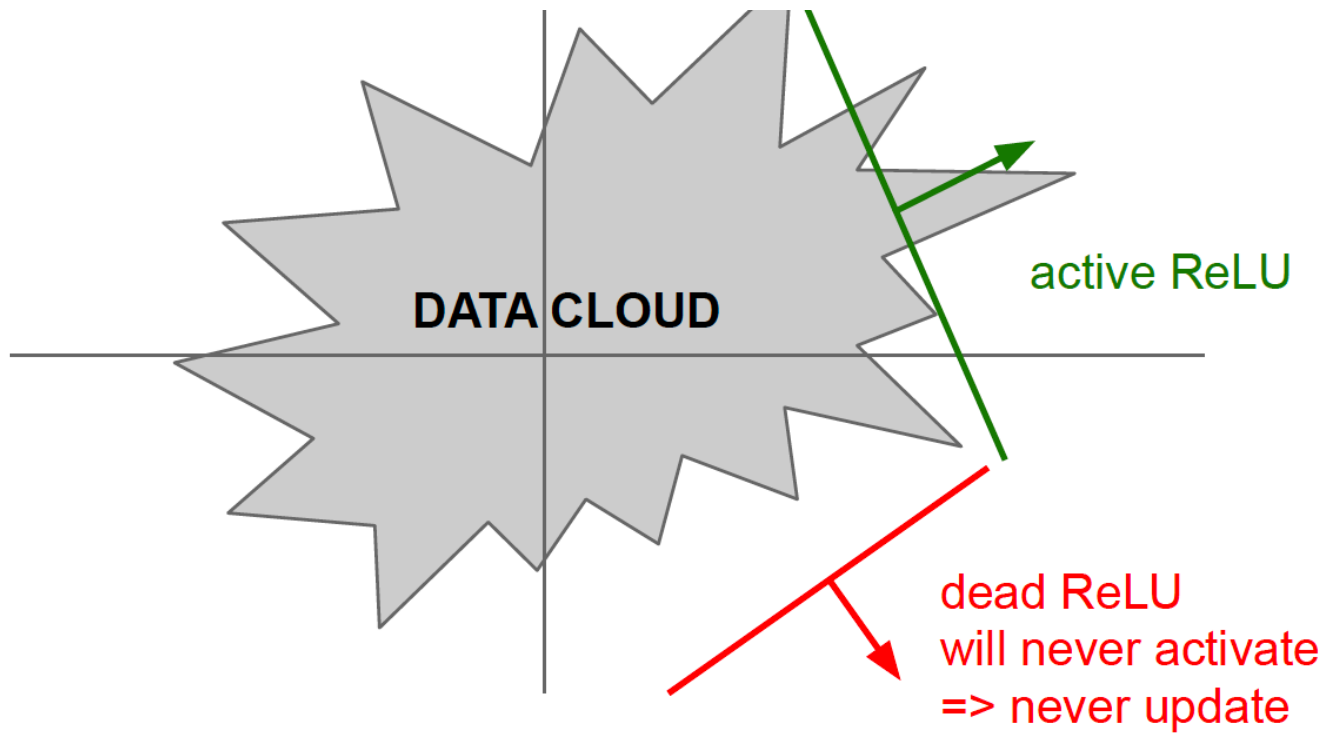
Problems

1. Not zero-centered output
2. Non-differentiable at zero; however, it is differentiable anywhere else.

[Krizhevsky et al., 2012]

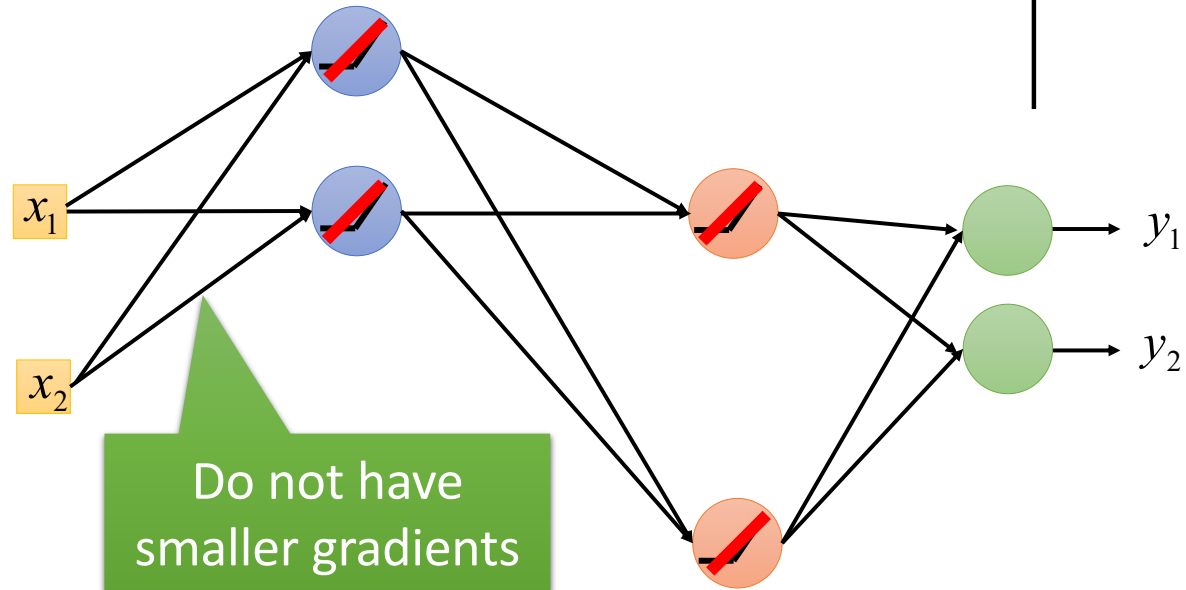
Rectified Linear Unit (ReLU)





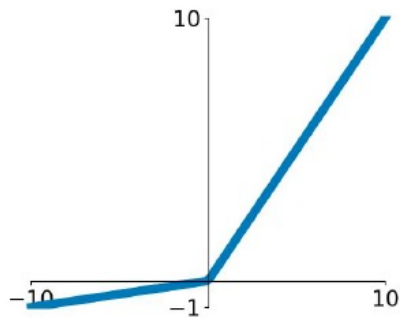
Rectified Linear Unit (ReLU)

A Thinner linear network



Activation Functions

[Mass et al., 2013]
[He et al., 2015]



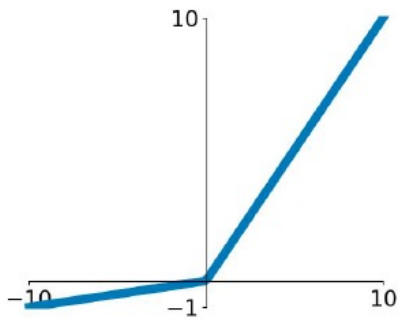
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Leaky ReLU

$$f(x) = \max(0.01x, x)$$

Activation Functions

[Mass et al., 2013]
[He et al., 2015]



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

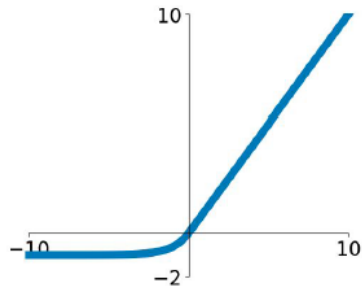
Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into α
(parameter)

Activation Functions

Exponential Linear Units (ELU)



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

(Alpha default = 1)

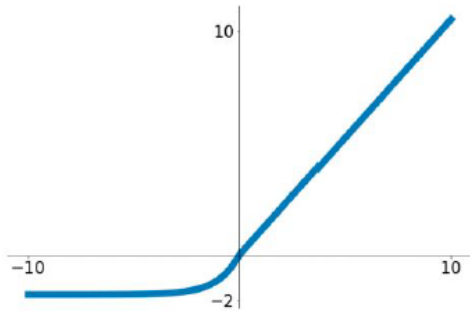
- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

- Computation requires $\exp()$

[Clevert et al., 2015]

Activation Functions

Scaled Exponential Linear Units (SELU)



$$f(x) = \begin{cases} \lambda x & \text{if } x > 0 \\ \lambda \alpha (e^x - 1) & \text{otherwise} \end{cases}$$

$$\alpha = 1.6733, \lambda = 1.0507$$

- Scaled version of ELU that works better for deep networks
- “Self-normalizing” property;
- Can train deep SELU networks without BatchNorm
 - (will discuss more later)

[Klambauer et al. ICLR 2017]

Maxout “Neuron”

- Does not have the basic form of dot product -> nonlinearity
- Generalizes ReLU and Leaky ReLU
- Linear Regime! Does not saturate! Does not die!

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

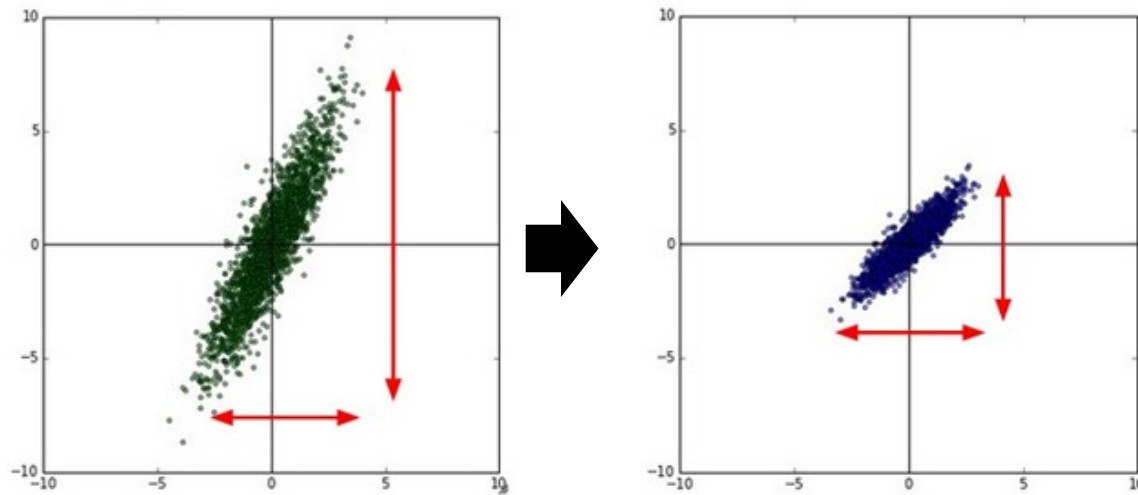
Problem: doubles the number of parameters/neuron :(

In Practice

- Use ReLU
- Don't use sigmoid or tanh

Data Preprocessing-Feature Scaling

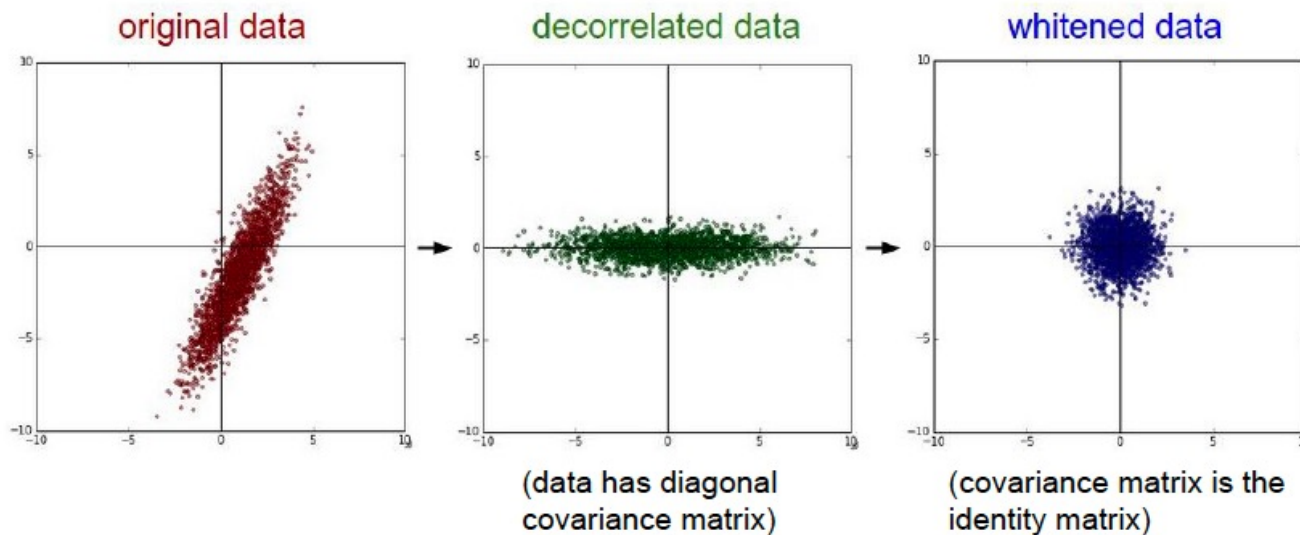
$$y = b + w_1x_1 + w_2x_2$$



Make different features have the same scaling

Data Preprocessing-Feature Scaling

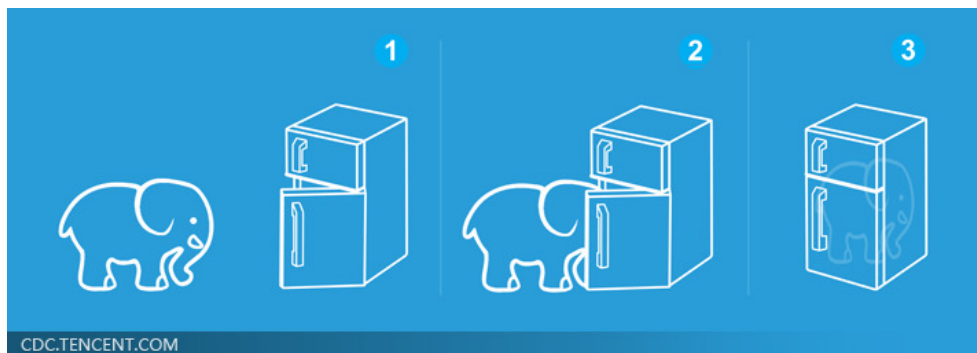
In practice, you may also see **PCA** and **Whitening** of the data



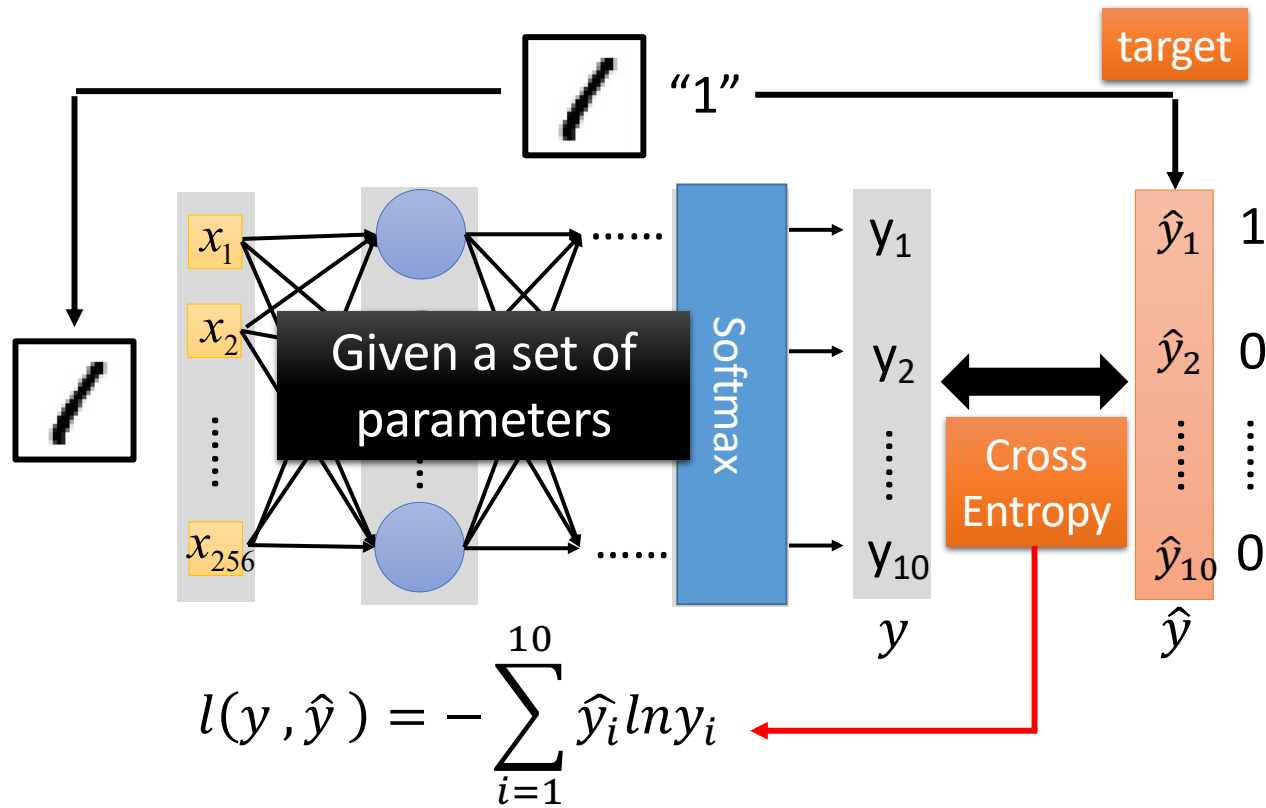
Three steps for deep learning



Deep Learning is so simple



Loss for an Example

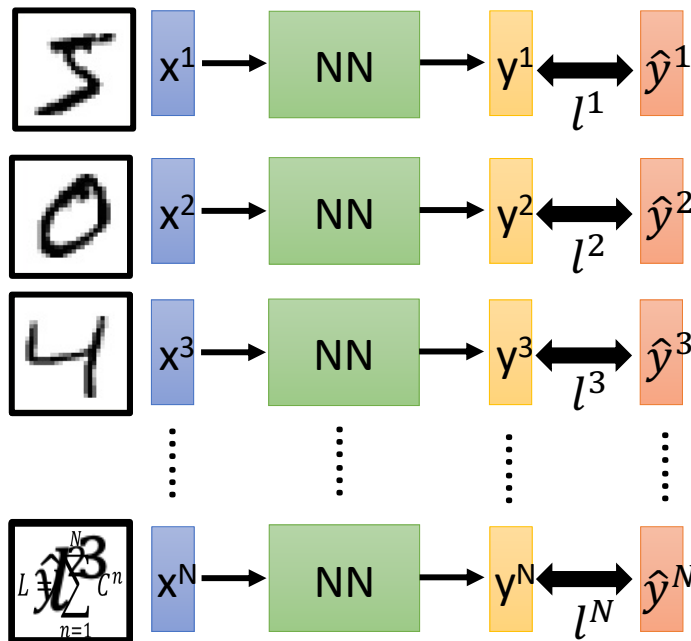


Total Loss

Total Loss:

$$L = \sum_{n=1}^N l^n$$

For all training data ...



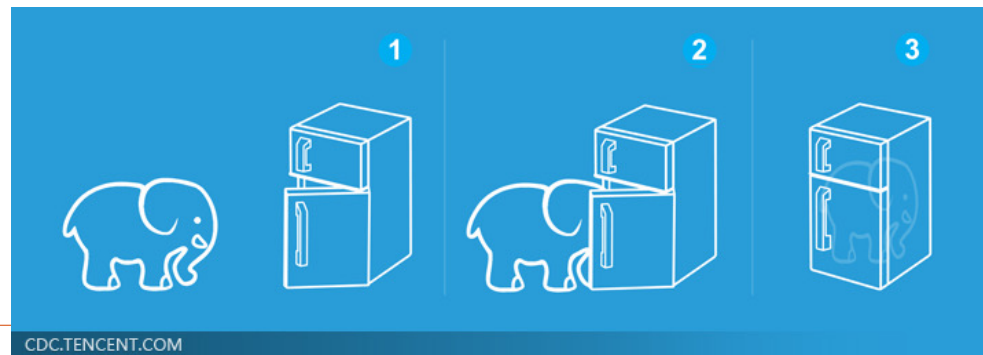
Find a function in function set that minimizes total loss L

Find the network parameters θ^* that minimize total loss L

Three Steps for Deep Learning



Deep Learning is so simple



Gradient Descent



Gradient Descent

- In step 3, we have to solve the following optimization problem:

$$\theta^* = \arg \min_{\theta} L(\theta) \quad L: \text{loss function} \quad \theta: \text{parameters}$$

Suppose that θ has two variables $\{\theta_1, \theta_2\}$

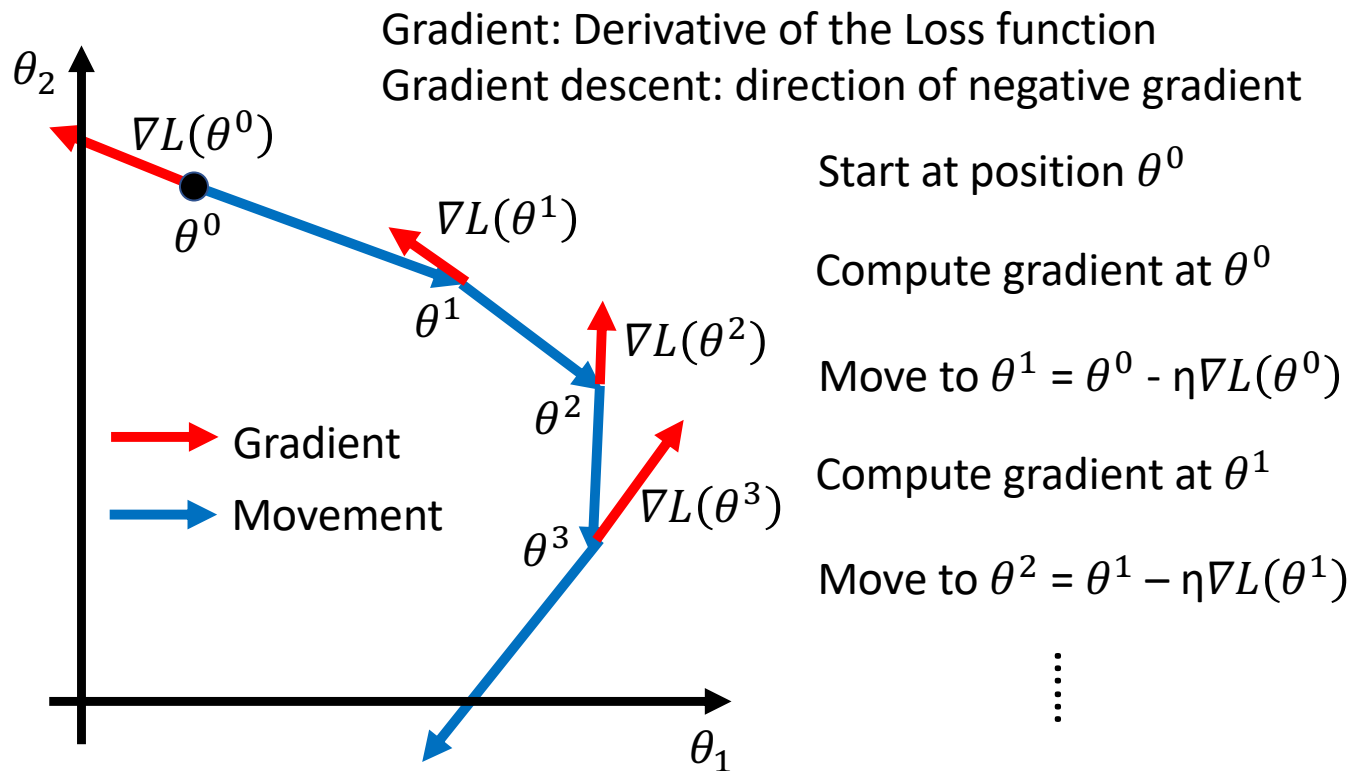
Randomly start at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta_1)/\partial \theta_1 \\ \partial L(\theta_2)/\partial \theta_2 \end{bmatrix}$$

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^0)/\partial \theta_1 \\ \partial L(\theta_2^0)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^1 = \theta^0 - \eta \nabla L(\theta^0)$$

$$\begin{bmatrix} \theta_1^2 \\ \theta_2^2 \end{bmatrix} = \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^1)/\partial \theta_1 \\ \partial L(\theta_2^1)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^2 = \theta^1 - \eta \nabla L(\theta^1)$$

Gradient Descent

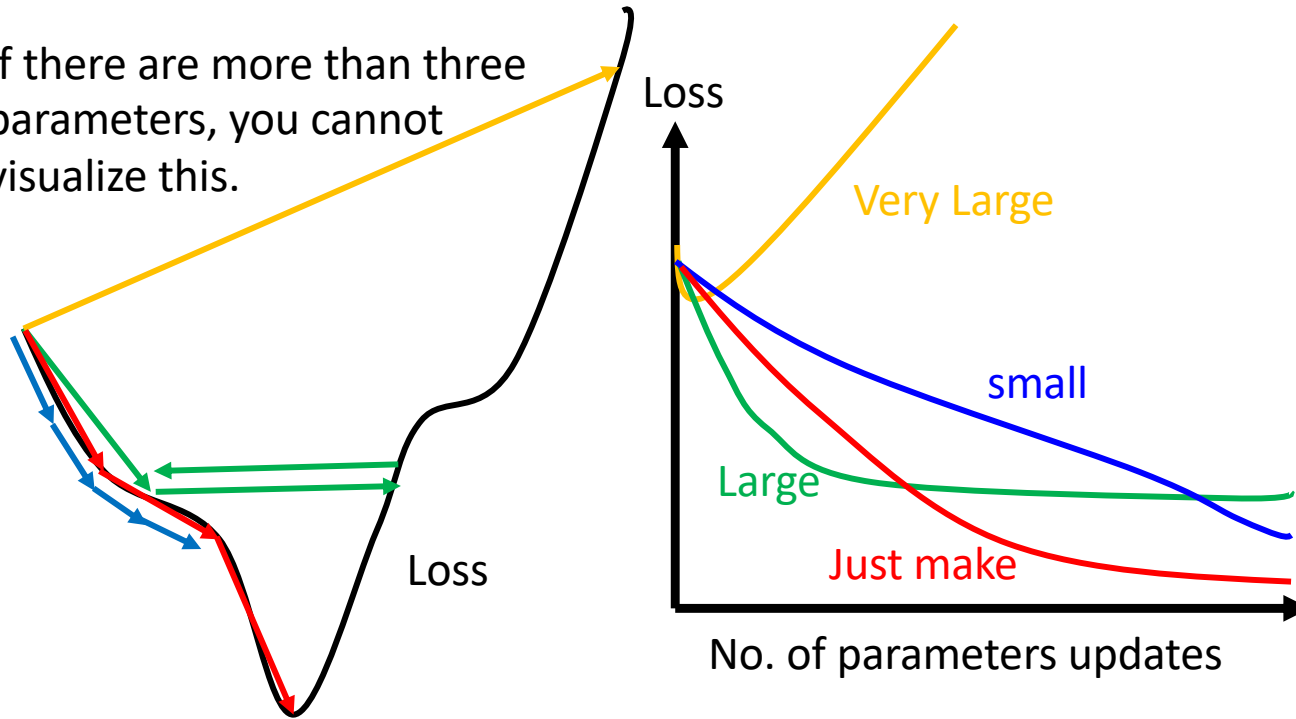


Tip 1: Tuning your learning rates

$$\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$$

Set the **learning rate η** carefully

If there are more than three parameters, you cannot visualize this.



But you can always visualize this.

Adaptive Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
 - At the beginning, we are far from the destination, so we use larger learning rate
 - After several epochs, we are close to the destination, so we reduce the learning rate
 - E.g. 1/t decay: $\eta^t = \eta / \sqrt{t + 1}$
- Learning rate cannot be one-size-fits-all
 - Giving different parameters different learning rates

Adaptive Learning Rates (Adagrad) $\eta^t = \frac{\eta}{\sqrt{t+1}}$ $g^t = \frac{\partial L(\theta^t)}{\partial w}$

- Divide the learning rate of each parameter by the *root mean square of its previous derivatives*

Vanilla Gradient descent (a.k.a. batch gradient descent)

$$w^{t+1} \leftarrow w^t - \eta^t g^t$$

w is one parameter

Adagrad

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

σ^t : root mean square of the previous derivatives of parameter w

$$\theta^{t+1} = \theta^t - \eta \nabla L(\theta^t)$$

Parameter dependent

Adagrad

σ^t : *root mean square* of the previous derivatives of parameter w

$$w^1 \leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0$$

$$\sigma^0 = \sqrt{(g^0)^2}$$

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1$$

$$\sigma^1 = \sqrt{\frac{1}{2} [(g^0)^2 + (g^1)^2]}$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2$$

$$\sigma^2 = \sqrt{\frac{1}{3} [(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

⋮

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$$\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$

Adagrad

- Divide the learning rate of each parameter by the **root mean square of its previous derivatives**

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$\eta^t = \frac{\eta}{\sqrt{t+1}}$ 1/t decay

$\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

Stochastic Gradient Descent



Stochastic Gradient Descent

$$L = \sum_n \left(\hat{y}^n - \left(b + \sum w_i x_i^n \right) \right)^2$$

Loss is the summation over all training examples

◆ Gradient Descent $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$

◆ Stochastic Gradient Descent **Faster!**

Pick an example x^n

Drawn uniform at random from $\{1, \dots, n\}$

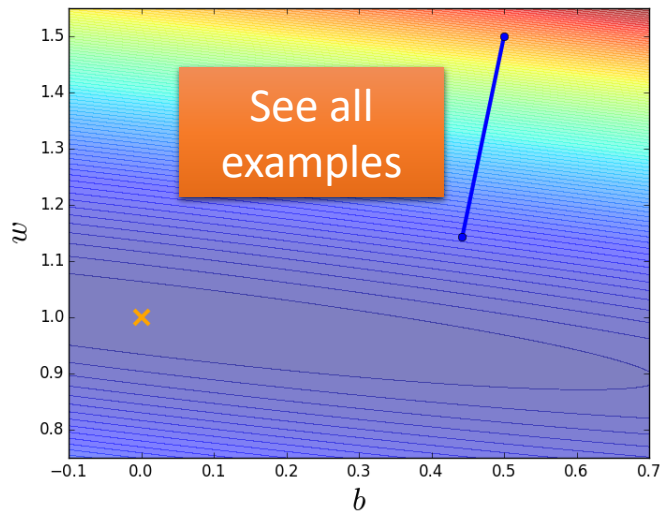
$$L^n = \left(\hat{y}^n - \left(b + \sum w_i x_i^n \right) \right)^2 \quad \theta^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$$

Loss for only **one example**, NO summing

SGD: An Illustration

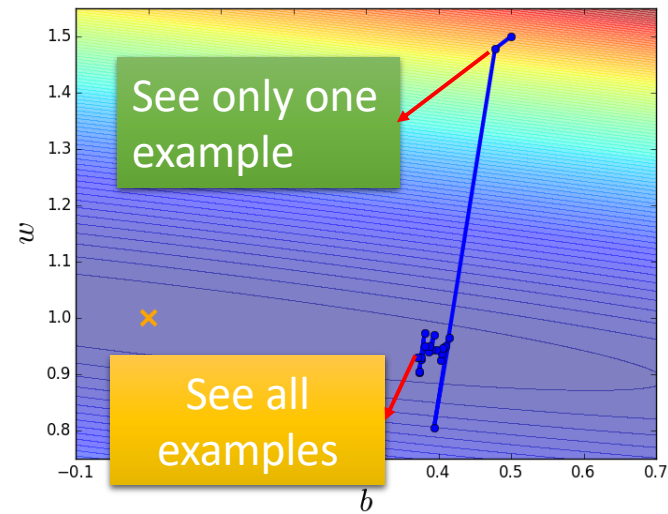
Gradient Descent

Update after seeing all examples



Stochastic Gradient Descent

Update for each example
If there are 20 examples,
20 times faster.

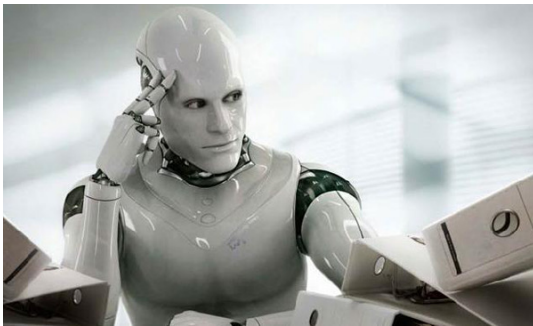


Gradient Descent

This is the “learning” of machines in deep learning

➔ Even alpha go using this approach.

People image



Actually

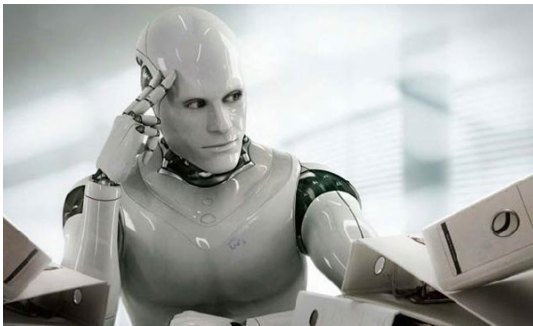


Gradient Descent

This is the “learning” of machines in deep learning

➔ Even alpha go using this approach.

People image



Actually



I hope you are not too disappointed :p

Backpropagation

- Backpropagation: an efficient way to compute $\partial L / \partial w$ in neural network



theano

Caffe

Microsoft
CNTK

PyTorch

Deep Learning library produced by Amazon

DSSTNE

mxnet

Any Questions?

bidong@syr.edu