# Introduction to Neural Networks

Philipp Seegerer and Seul-Ki Yeom
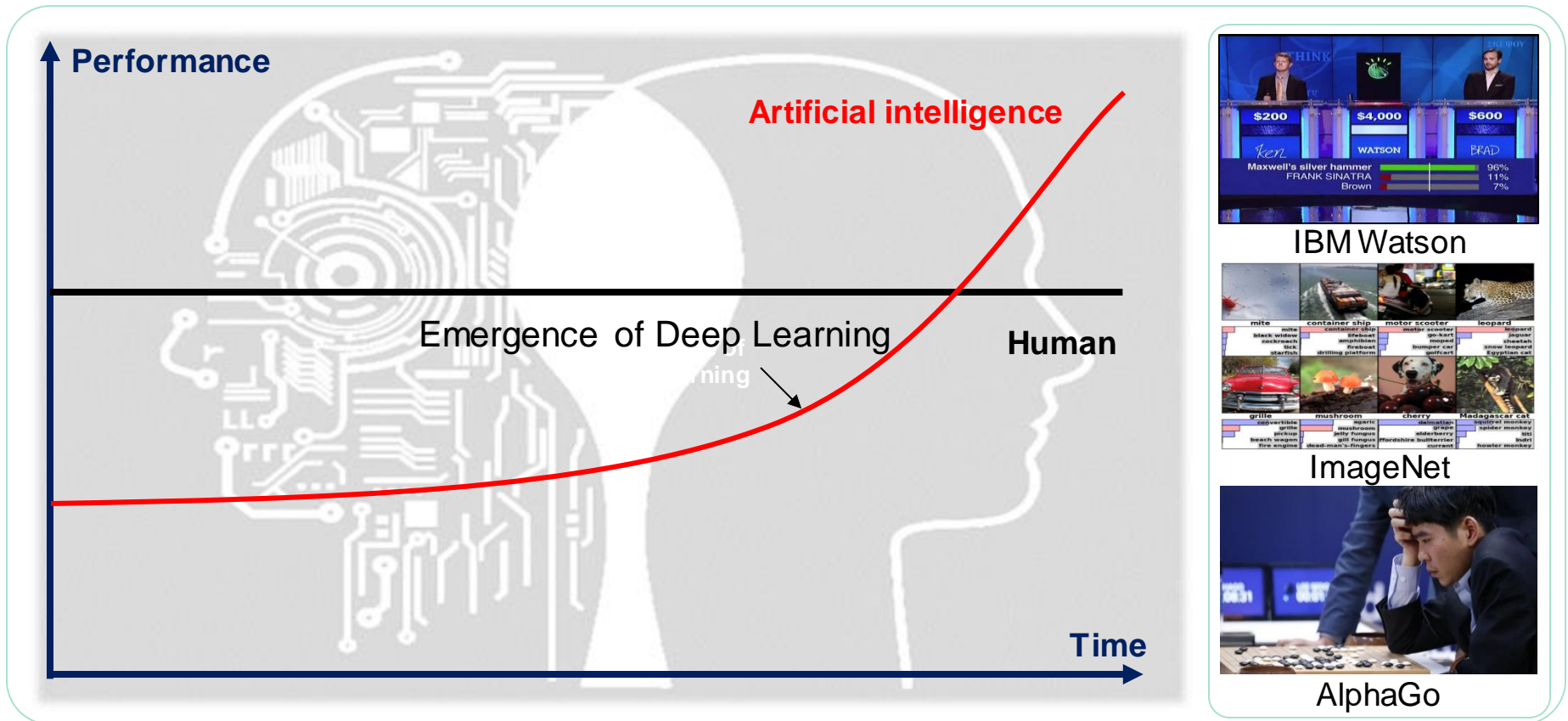
Technische Universität Berlin - Machine Learning Group

2019 Beginners Workshop Machine Learning

# Contents

- History of NN

- Rosenblatt's Perceptron

  - Cost (= Loss) function

  - Gradient descent algorithm

- Multi-Layer Perceptron (MLP)
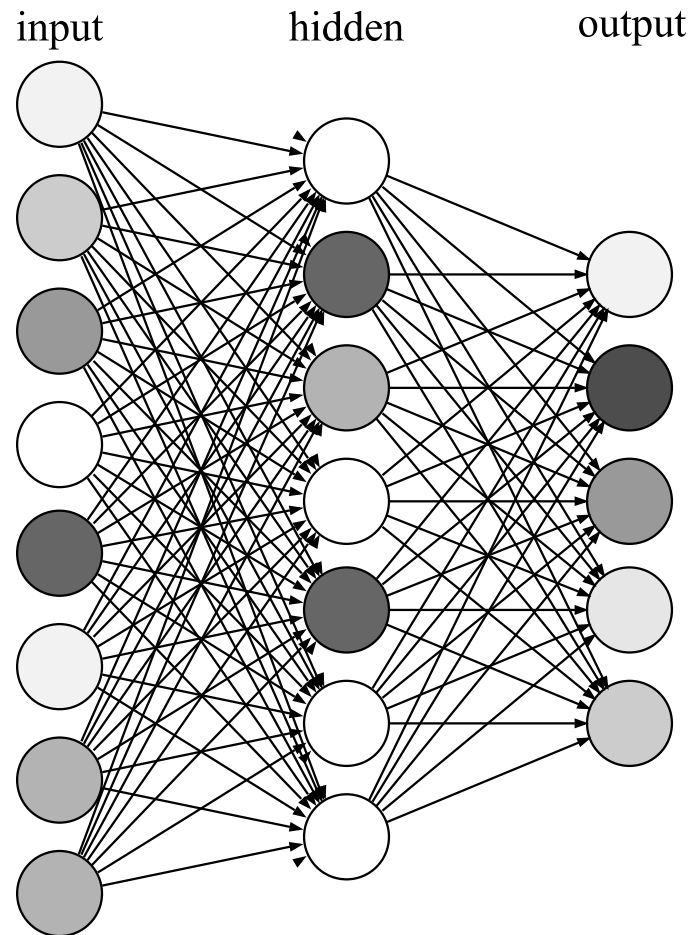
  - Activation function

  - Backpropagation

# Emergence of Deep Learning



- The technology of AI has been remarkably developed and now has surpassed human intelligence
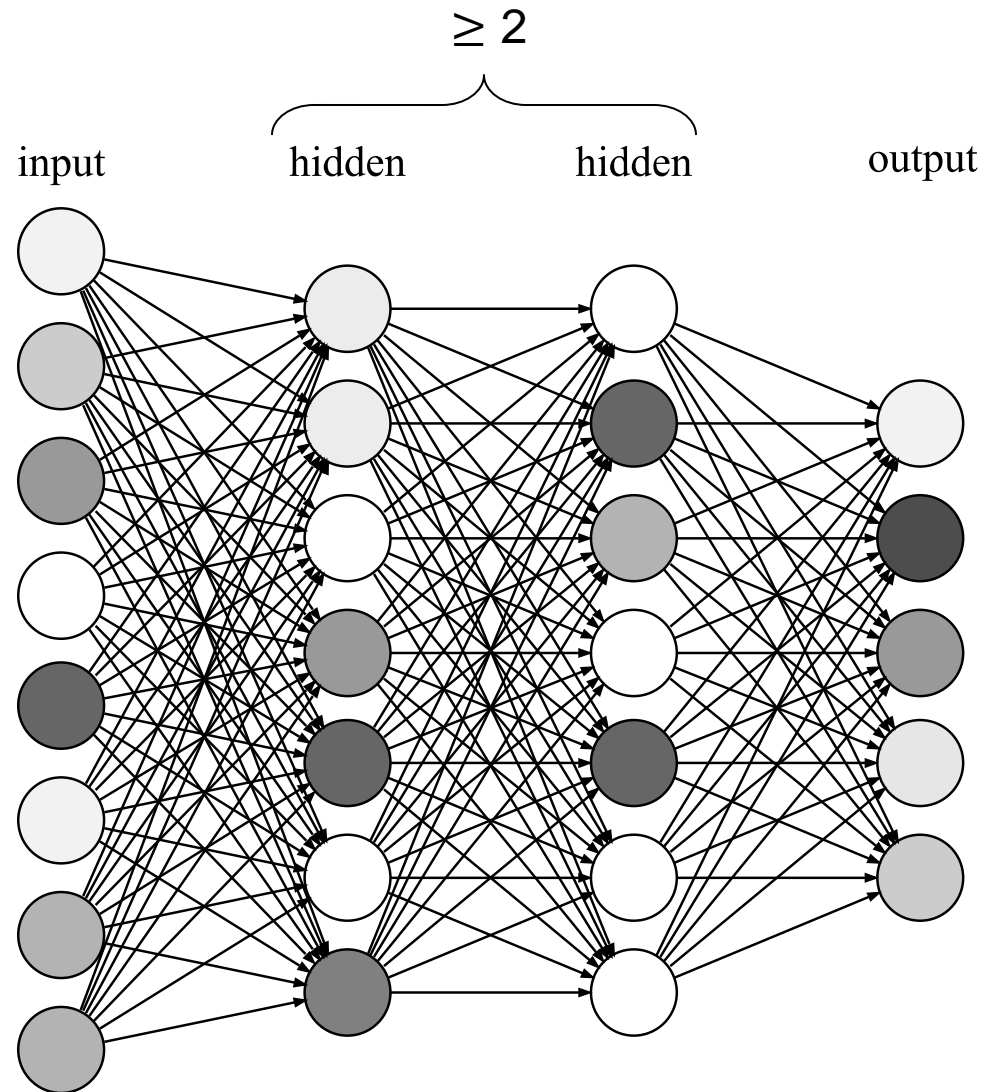
# (Artificial) Neural Network

- Deep learning is one of the ML research tool which uses NN
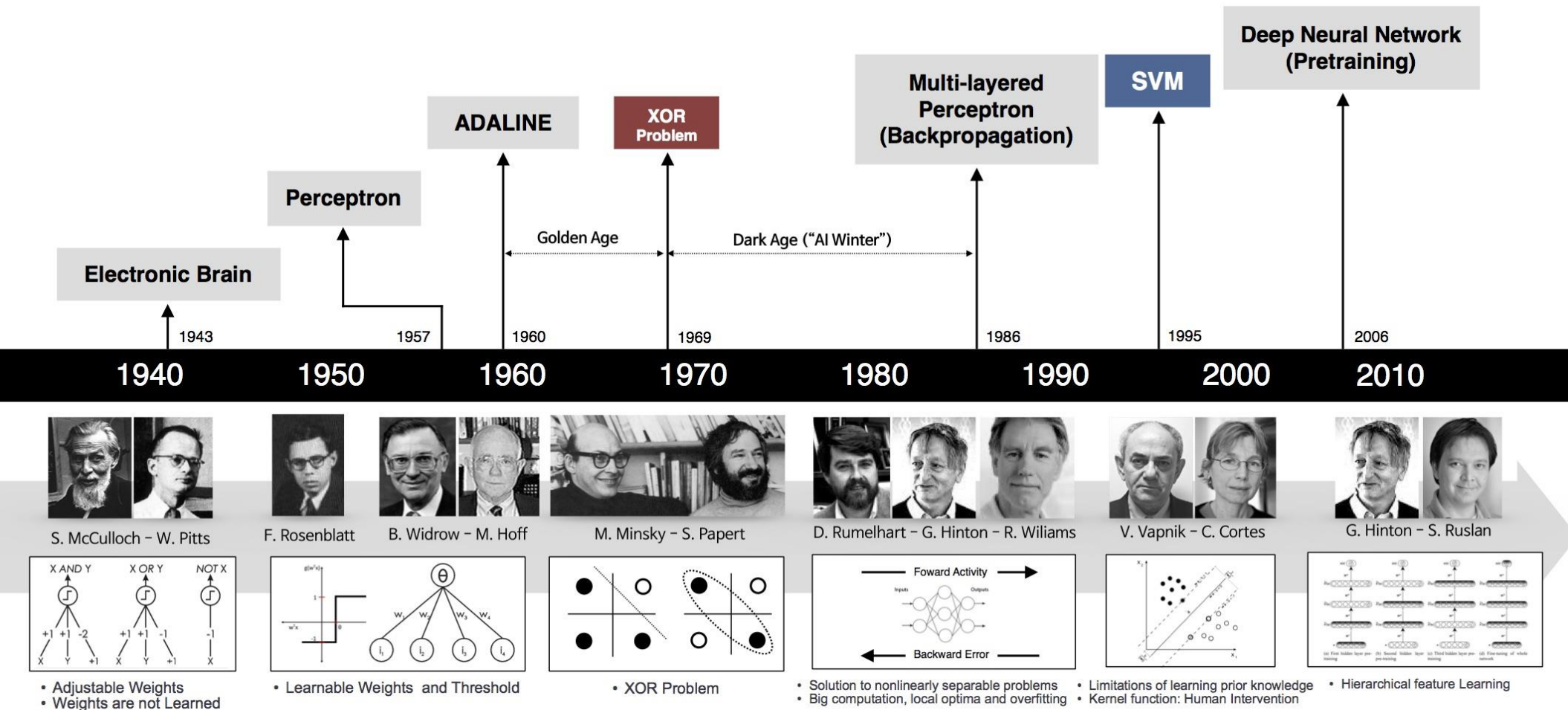
- Pieces of math !

input      hidden      output
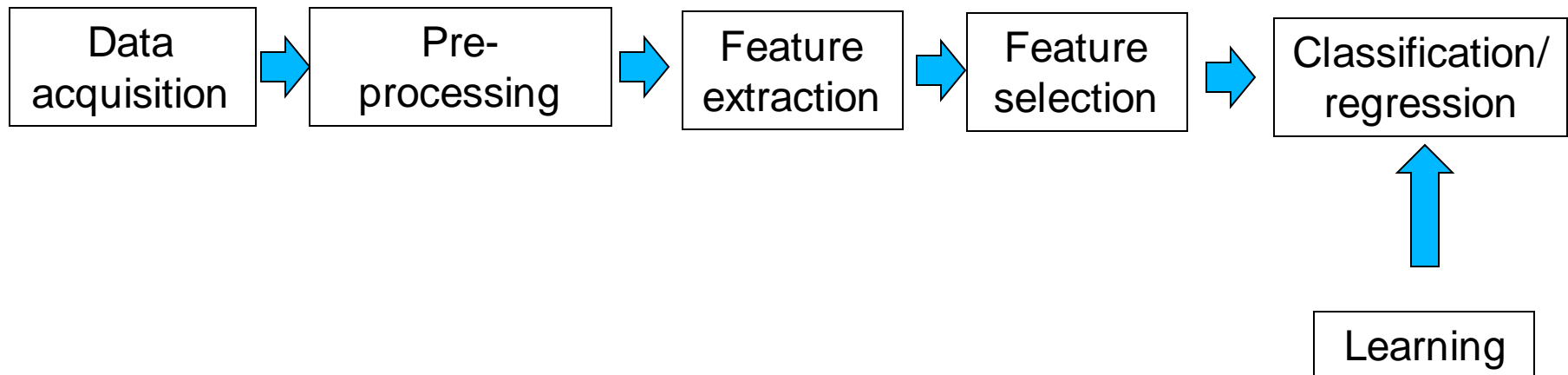
Multi-layer Perceptron

# Deep Neural Network

- # of hidden layer $\geq$ 2

# Milestones in the Development of Neural Networks

https://beamandrew.github.io/deeplearning/2017/02/23/deep_learning_101_part1.html

# Core Idea: Feature Learning

Classical pattern recognition pipeline

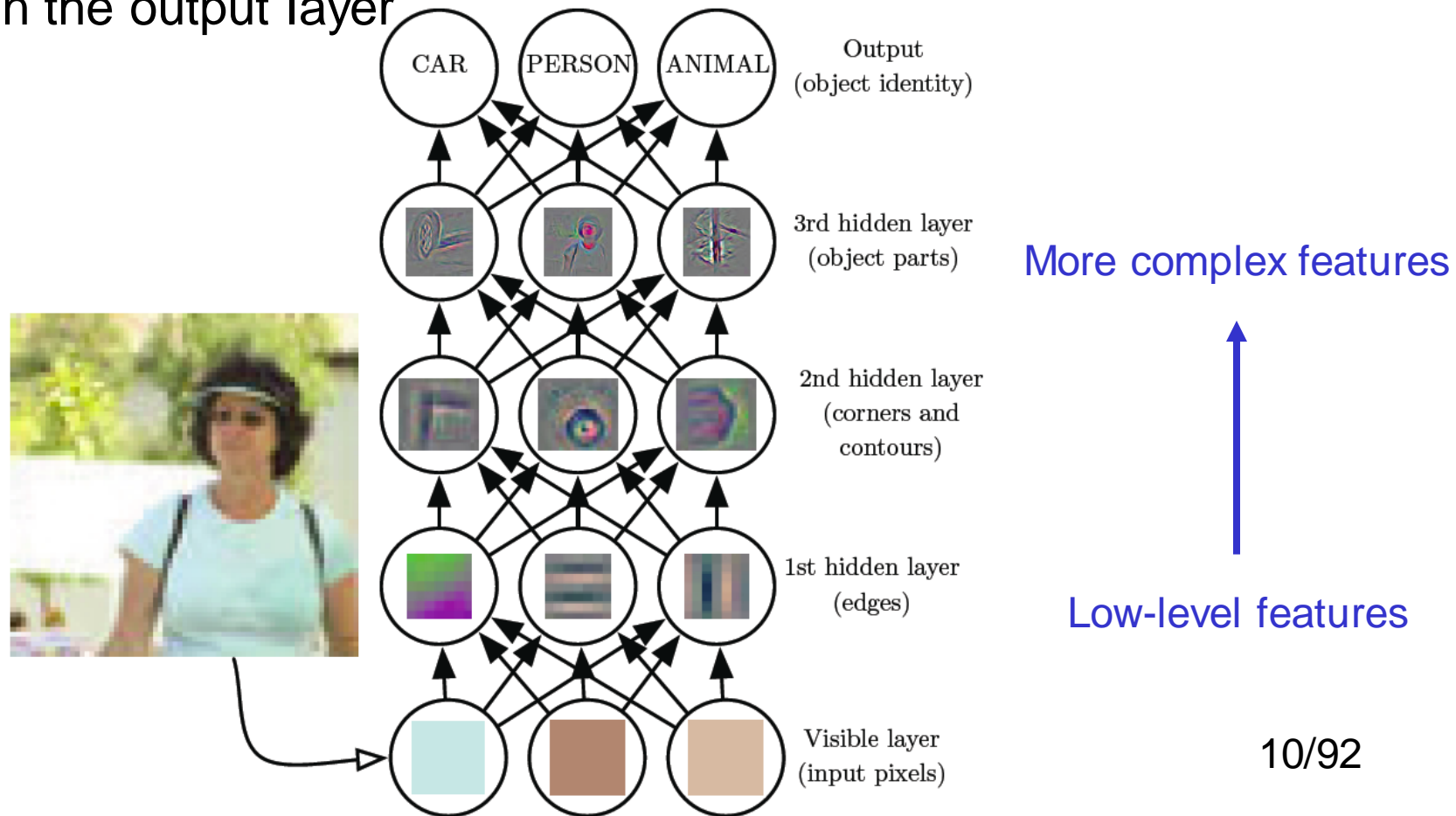| Data acquisition | → | Pre-processing | → | Feature extraction | → | Feature selection | → | Classification/regression |
|---|---|---|---|---|---|---|---|---|

Learning

# Core Idea: Feature Learning
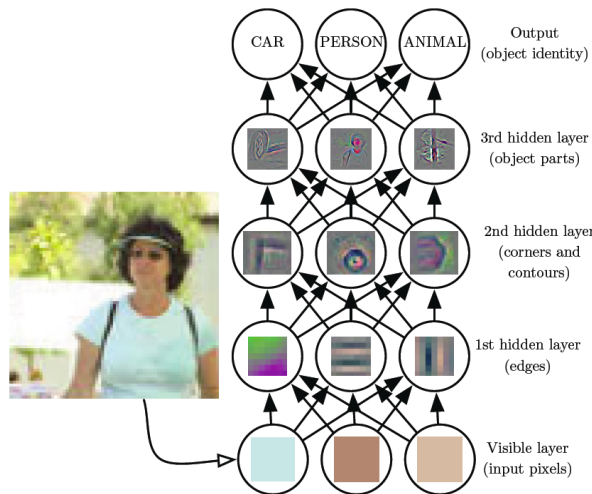
Neural networks pipeline

# Deep Neural Network in action

- Learning representations with increasing level of abstraction

- By passing it with several layers hierarchically, we can classify the images in the output layer



More complex features

Low-level features

Goodfellow, 2016

# Deep Neural Network in action

- Image recognition

  - pixel → edge → Texton → motif → part → object

- Text

  - Character → word → word group → clause → sentence → story

- Speech

  - sample → spectral band → sound → … → phone → phoneme → word
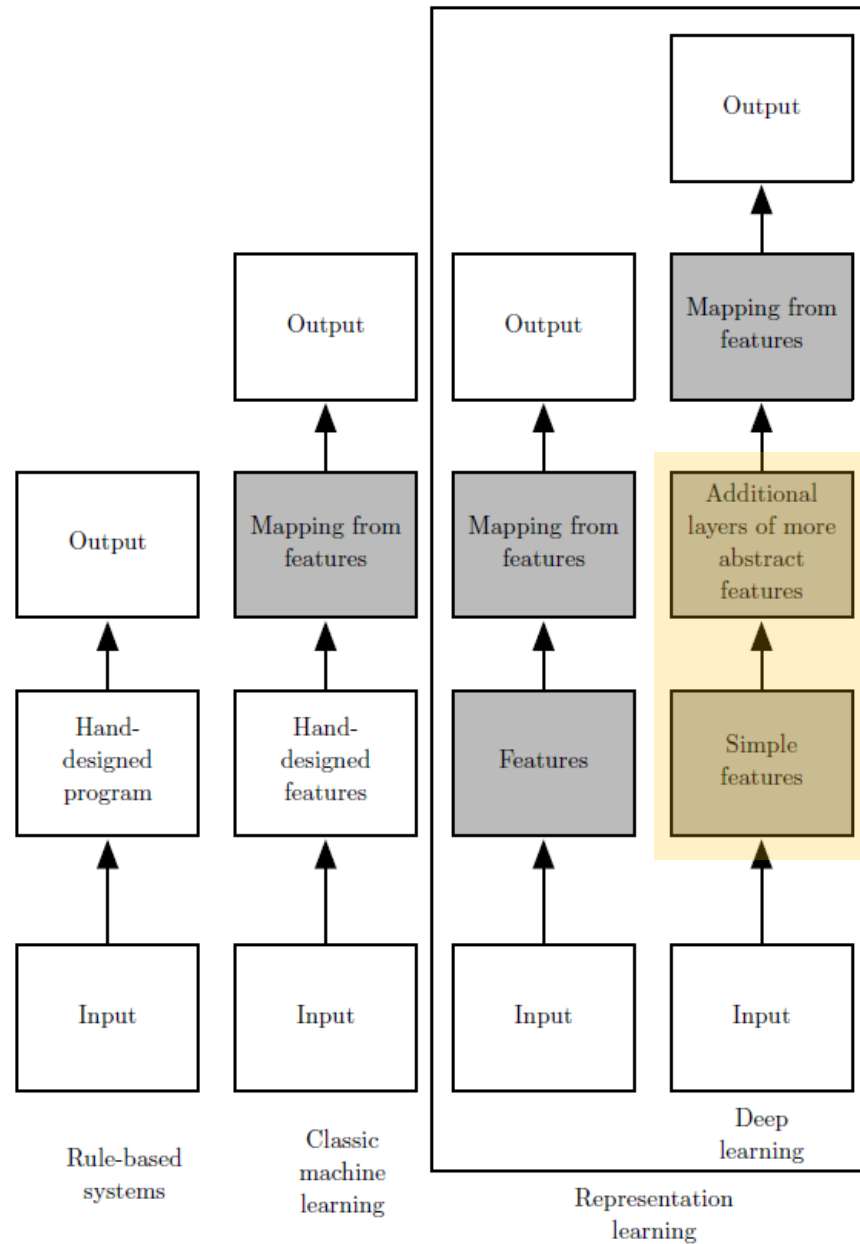


More complex features

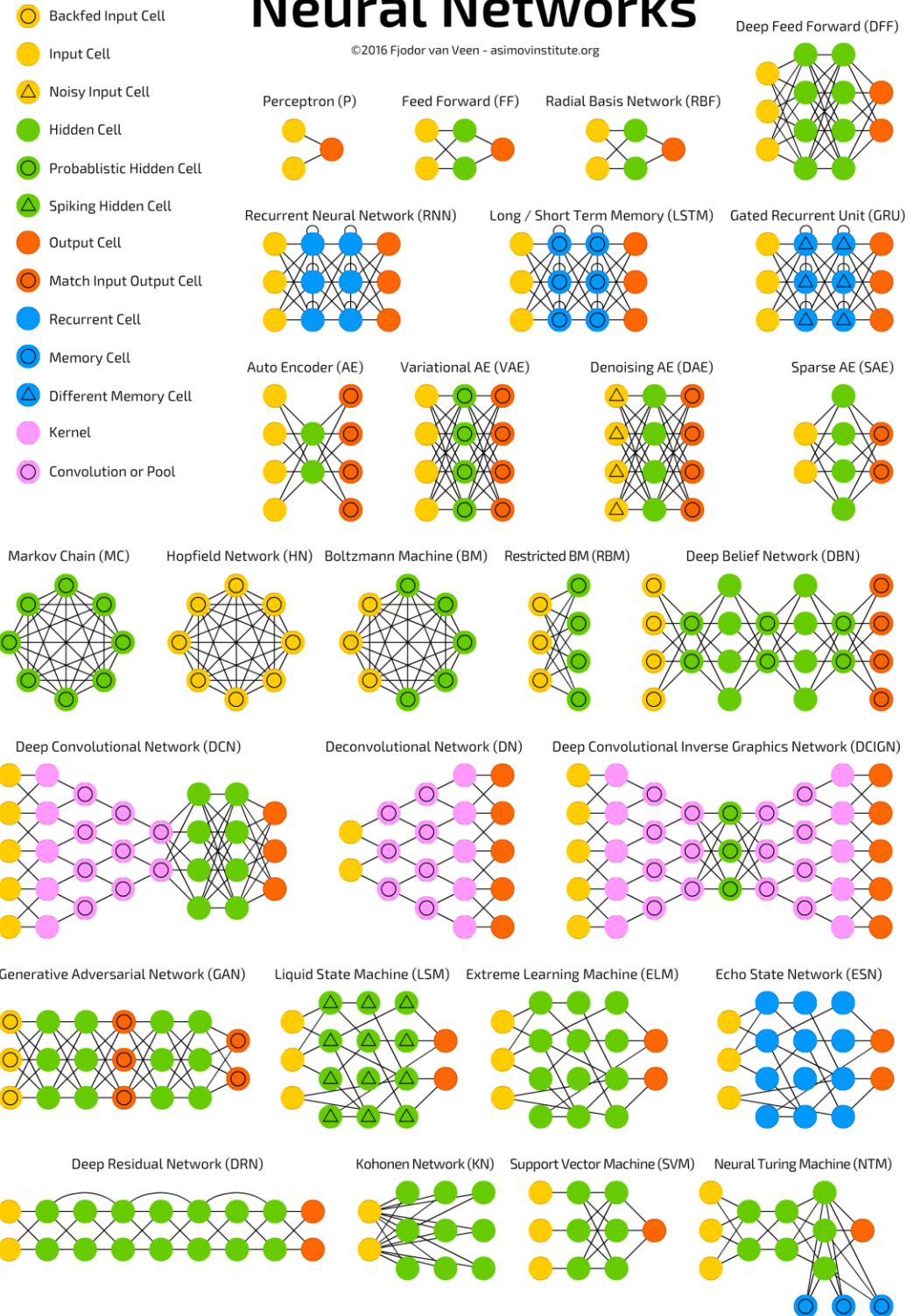Low-level features

# Learning Multiple Components



Automatic feature extraction

Goodfellow, 2016

# Neural Network Zoo



A mostly complete chart of
# Neural Networks
©2016 Fjodor van Veen - asimovinstitute.org

**Legend:**
- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool

Deep Feed Forward (DFF)

Perceptron (P)  Feed Forward (FF)  Radial Basis Network (RBF)

Recurrent Neural Network (RNN)  Long / Short Term Memory (LSTM)  Gated Recurrent Unit (GRU)

Auto Encoder (AE)  Variational AE (VAE)  Denoising AE (DAE)  Sparse AE (SAE)

Markov Chain (MC)  Hopfield Network (HN)  Boltzmann Machine (BM)  Restricted BM (RBM)  Deep Belief Network (DBN)

Deep Convolutional Network (DCN)  Deconvolutional Network (DN)  Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)  Liquid State Machine (LSM)  Extreme Learning Machine (ELM)  Echo State Network (ESN)

Deep Residual Network (DRN)  Kohonen Network (KN)  Support Vector Machine (SVM)  Neural Turing Machine (NTM)

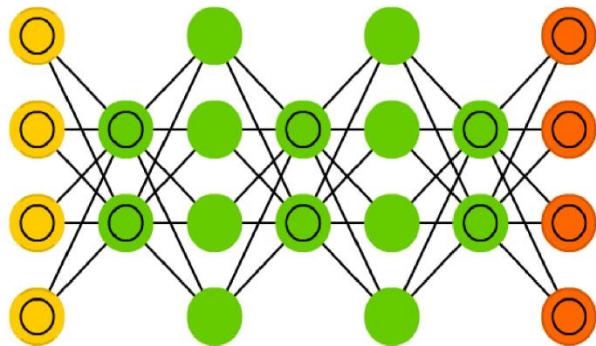http://www.asimovinstitute.org/neural-network-zoo/

# Deep Neural Networks: Applications and Models



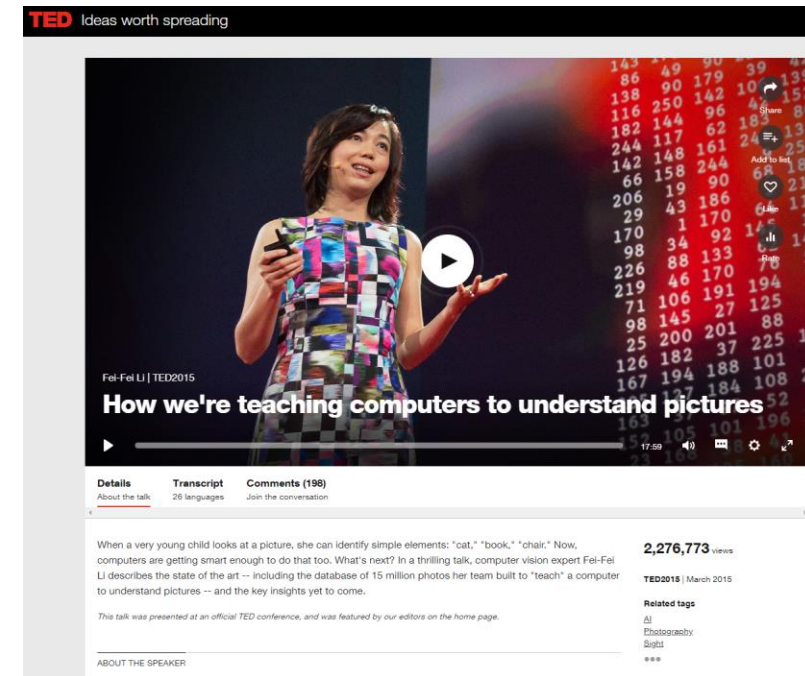| Internet & Cloud | Medicine & Biology | Media & Entertainment | Security & Defense | Autonomous Machines |
|---|---|---|---|---|
| Image classification Speech recognition Language Translation Language processing Sentiment analysis Recommendation | Cancer cell detection Diabetic grading Drug discovery | Video captioning Video search Real time translation | Face detection Video surveillance Satellite imagery | Pedestrian Detection Lane Tracking Traffic sign recognition |



Deep Belief Network



Convolutional Neural Network



Recurrent Neural Network

# References on Neural Networks

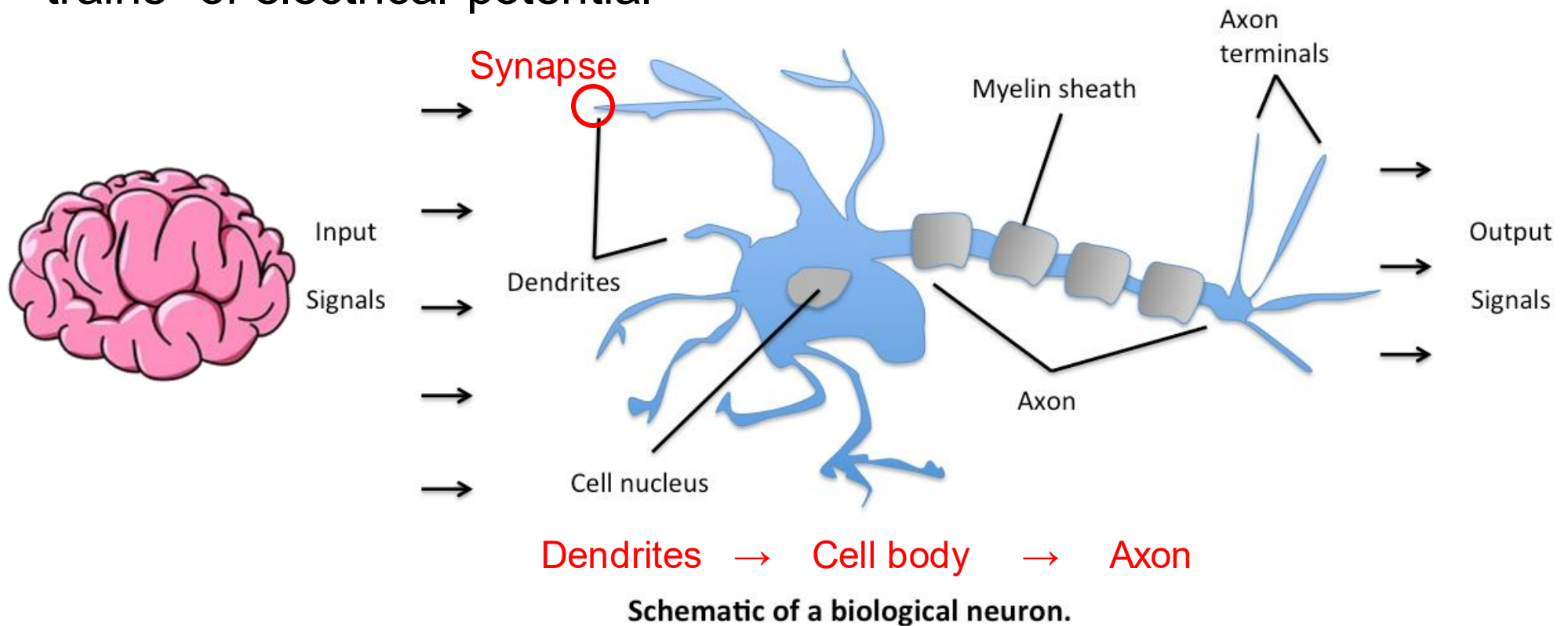- Andrew Ng's and other ML tutorials

  - https://class.coursera.org/ml-003/lecture

  - http://www.holehouse.org/mlclass/

  - http://deeplearning.stanford.edu/tutorial/

- Stanford lecture videos

  - CS231n: Convolutional Neural Networks for Visual Recognition (http://cs231n.stanford.edu/syllabus.html)

  - CS224d: Deep Learning for Natural Language Processing (http://cs224d.stanford.edu/syllabus.html)

- Ted by Fei-Fei Li (How we're teaching computers to understand pictures (https://www.ted.com/talks/fei_fei_li_how_we_re_teaching_computers_to_understand_pictures)
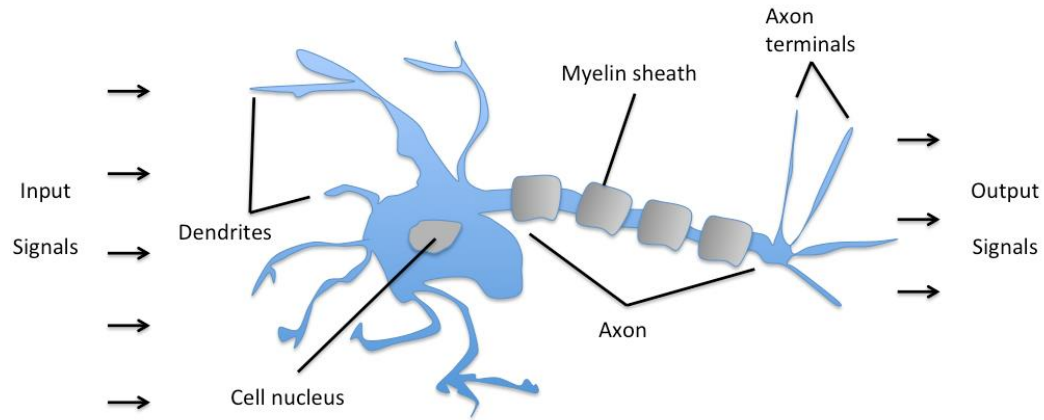
# Perceptron
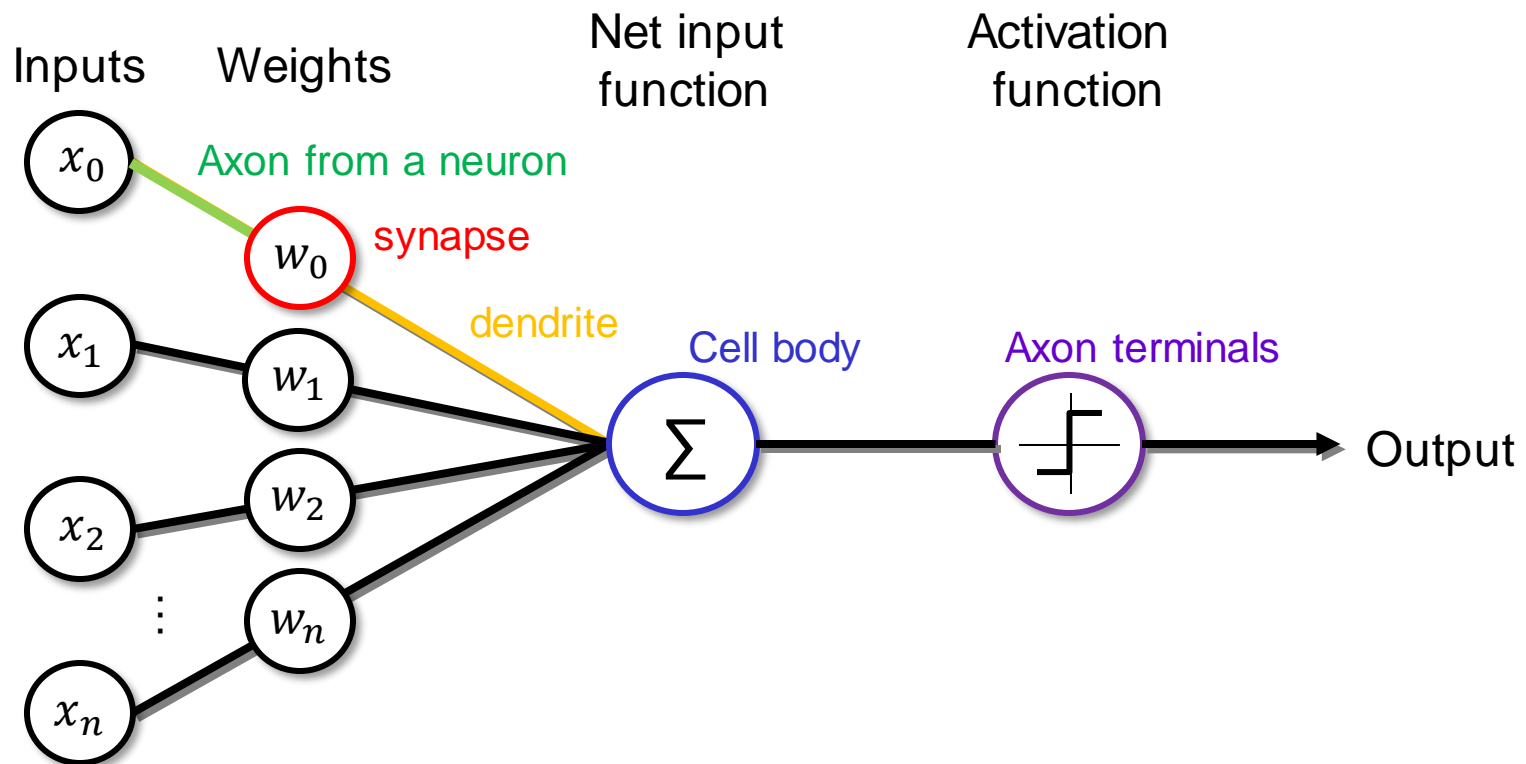
# Neuronal Activity in the Brain

$10^{11}$ neurons of > 20 types, $10^{14}$ synapses with very complex connections, 1ms–10ms cycle time Signals are noisy "spike trains" of electrical potential
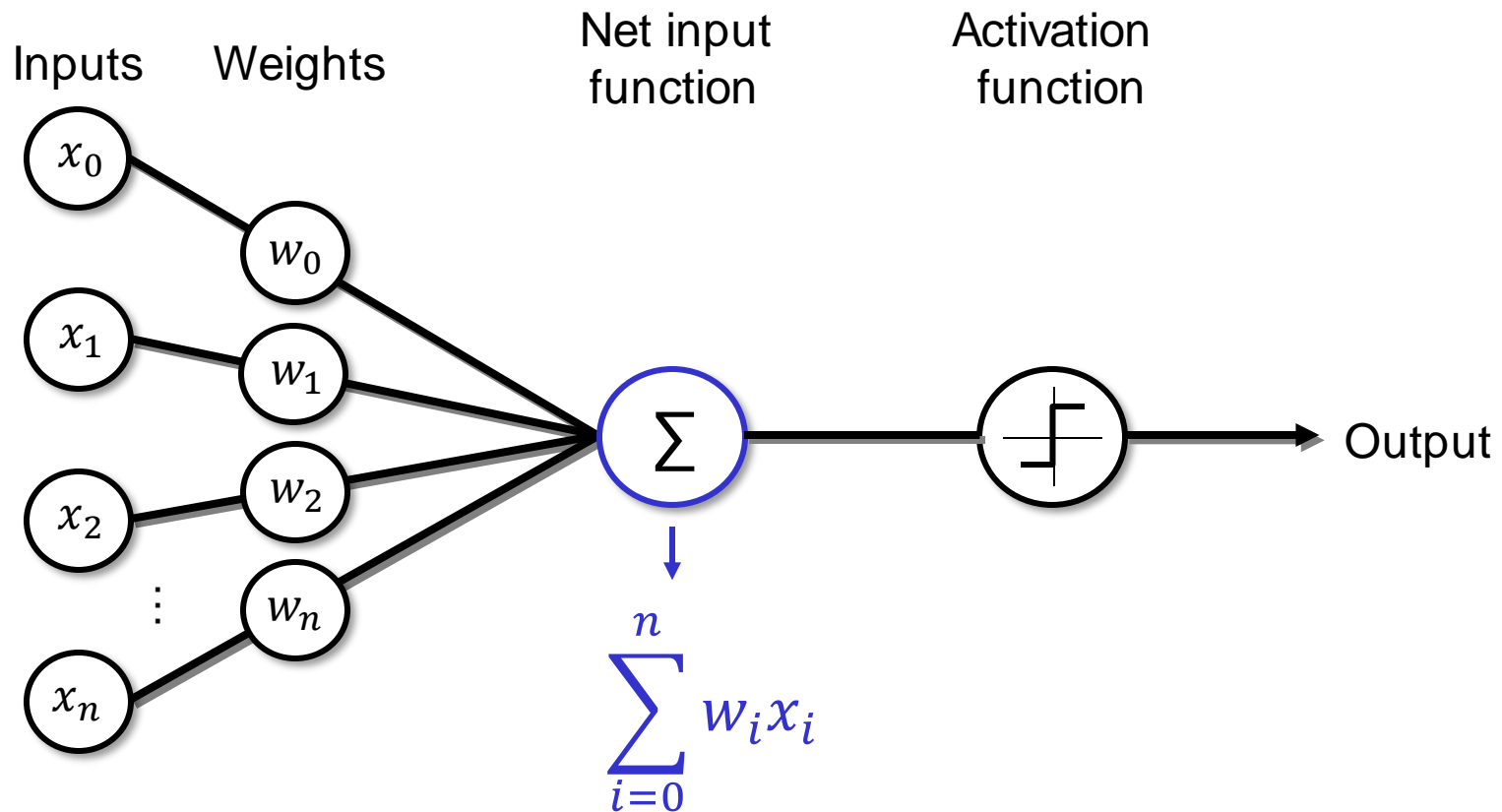


Dendrites  →  Cell body  →  Axon

**Schematic of a biological neuron.**

# Rosenblatt's Perceptron



Schematic of a biological neuron.

Inputs    Weights    Net input function    Activation function

$x_0$

Axon from a neuron

$w_0$    synapse

$x_1$    $w_1$    dendrite    Cell body    Axon terminals

$x_2$    $w_2$    $\Sigma$    Output

$\vdots$    $w_n$

$x_n$

Schematic of Rosenblatt's perceptron

18/92

# Rosenblatt's Perceptron: Cell Body

# Rosenblatt's Perceptron: Activation Function



$$\hat{y} = \begin{cases} +1 & if \sum_{i=0}^{n} w_i x_i \geq 0 \\ -1 & otherwise. \end{cases}$$
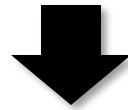
# Perceptron



Nerve cell

Schematic of a biological neuron.

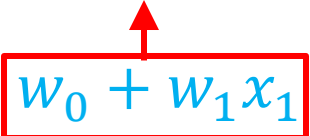**Linear Threshold Unit**

**(simple) Perceptron** $\hat{y}(x_1, \ldots, x_n) = \begin{cases} +1 & if\ \mathbf{w}^T\mathbf{x} \geq 0 \\ -1 & otherwise. \end{cases}$

# Perceptron

$$y = ax + b \quad \textit{Linear equation}$$

$$\hat{y}(x_1, \ldots, x_n) = \begin{cases} +1 & \textit{if } \boxed{w_0 + w_1 x_1} + \cdots + w_n x_n \geq 0 \\ -1 & \textit{otherwise}. \end{cases}$$
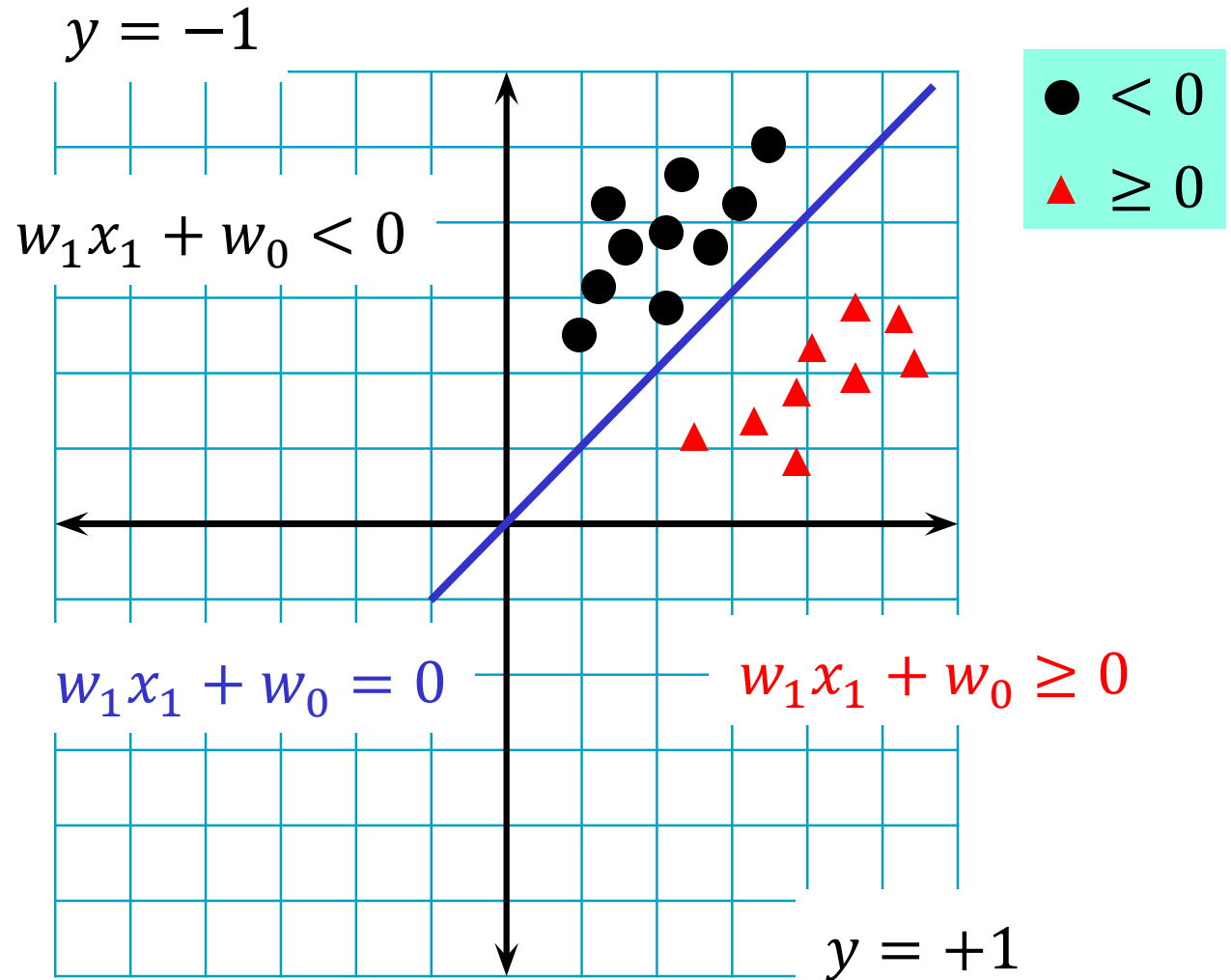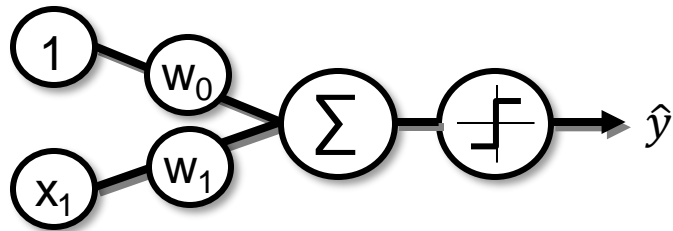
$$\hat{y}(x_1, \ldots, x_n) = \begin{cases} +1 & \textit{if } \boldsymbol{w}^T \mathbf{x} \geq 0 \\ -1 & \textit{otherwise}. \end{cases}$$
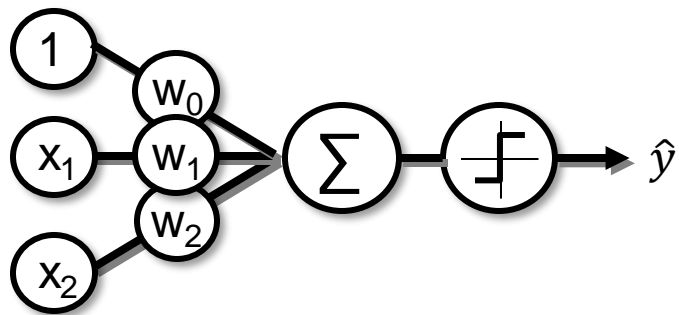
# Perceptron on 2-D coordinate

In case if $i = 1$



$y = -1$

$w_1 x_1 + w_0 < 0$

● $< 0$

▲ $\geq 0$

$w_1 x_1 + w_0 = 0$

$w_1 x_1 + w_0 \geq 0$

$y = +1$

# Perceptron on 2-D coordinate

In case if $i = 2$



$y = -1$

$\mathbf{w}^T\mathbf{x} < 0$
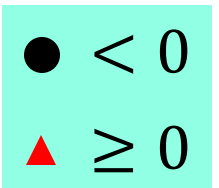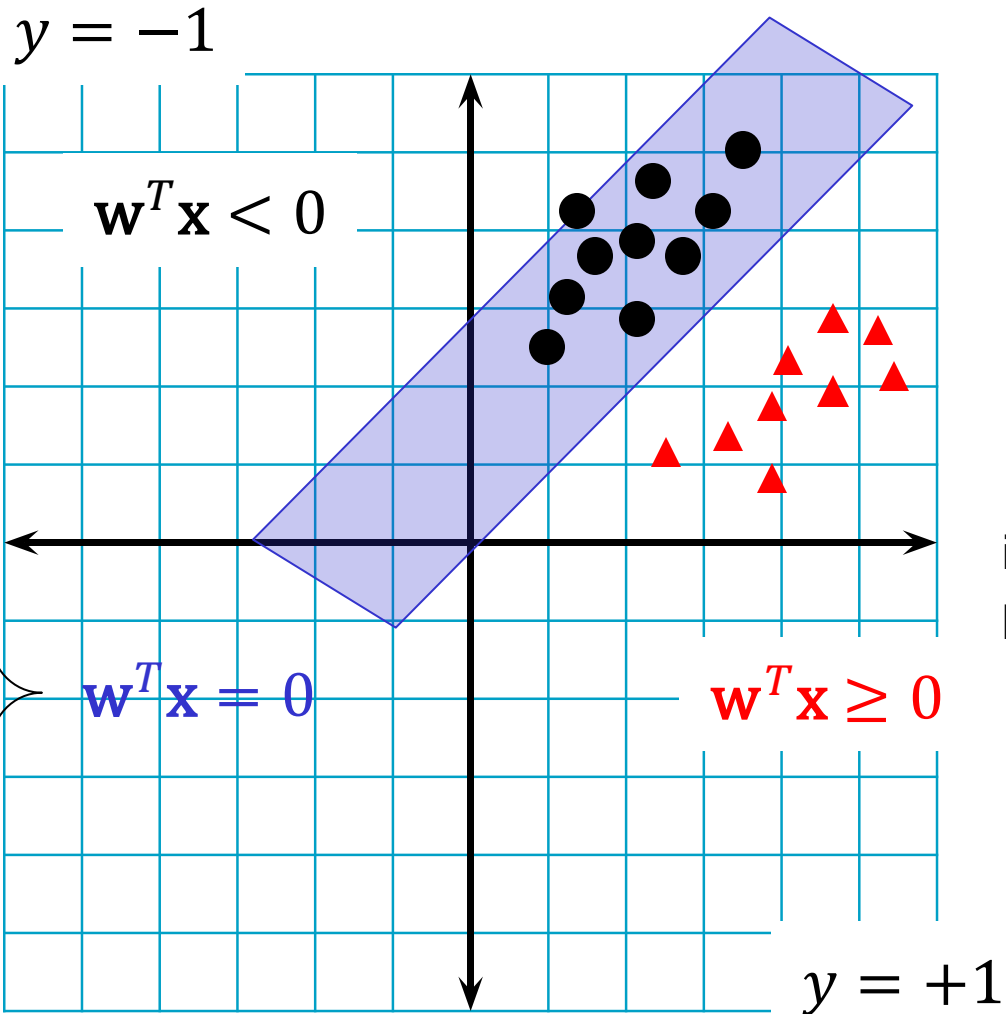
$$\bullet \quad < 0$$
$$\blacktriangle \quad \geq 0$$

$w_1 x_1 + w_0 x_0 = 0$

$(w_0 \ w_1)\begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = 0$

$\mathbf{w}^T = (w_0 \ w_1)$

$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix}$

$\mathbf{w}^T\mathbf{x} = 0$

$\mathbf{w}^T\mathbf{x} \geq 0$

if $i > 2$, hyper-plane

$y = +1$

# Learning on Perceptron

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

$$= w_0 x_0 + w_1 x_1 + \cdots + w_n x_n$$

Delta rule of Rosenblatt's perceptron

1. Initialize the weights

2. For each training sample $\left(\mathbf{x}^{(i)}, y^{(i)}\right)$,

   1. Compute the output $\hat{y}$

   2. Update the weight using $\Delta w_i$

3. Repeat 2nd step

# Delta Update (=Training) Rule

1. Initialize the weights

2. For each training sample $\left(\mathbf{x}^{(1)}, y^{(1)}\right), \ldots, \left(\mathbf{x}^{(d)}, y^{(d)}\right) \in D$, perform the following:

   1. Compute the prediction output $\hat{y} = \mathbf{w}^T \mathbf{x}$

   2. Update the weight $w_i^{new} \leftarrow w_i^{old} + \Delta w_i$

$$( \Delta w_i = -\eta \times \sum_{d \in D} \left(y^{(d)} - \hat{y}^{(d)}\right) \times x_i )$$

3. Repeat until the error is less than threshold

$$\Delta w_i = -\eta \times \sum_{d \in D} \left(y^{(d)} - \hat{y}^{(d)}\right) \times x_i$$

$y^{(d)}$: Target (correct) output for sample $d$ (0 or 1)
$\hat{y}^{(d)}$: Perceptron output (continuous value)
$\eta$: Learning rate (range [0 1])

# Perceptron: Cost(= Loss) Function

Loss: $E(\mathbf{w}) \equiv \dfrac{1}{2} \displaystyle\sum_{d \in D} \underline{\left(y^{(d)} - \hat{y}^{(d)}\right)}^2$

$y^{(d)}$: target (0 or 1)
$\hat{y}^{(d)}$: output
$\hat{y} = w_0 + w_1 x_1 + \cdots + w_1 x_1$

Difference between target value $y^{(d)}$ and output value $\hat{y}^{(d)}$ for training sample $d$

$$\underset{\mathbf{w}}{\text{minimize}}\ E(\mathbf{w})$$

Our objective is to find **w** which can minimize cost function

# How cost(W) looks like?

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} \left( y^{(d)} - \hat{y}^{(d)} \right)^2 \implies E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} (y^{(d)} - \mathbf{w}^T \mathbf{x}^{(d)})^2$$

| $\mathbf{x}$ | $y$ |
|:---:|:---:|
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |

- $w = 1, E(\mathbf{w})$ ?

# How cost(W) looks like?

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} (y^{(d)} - \mathbf{w}^T \mathbf{x}^{(d)})^2$$

| x | y |
|---|---|
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |

- $w = 1, E(\mathbf{w}) = 0$

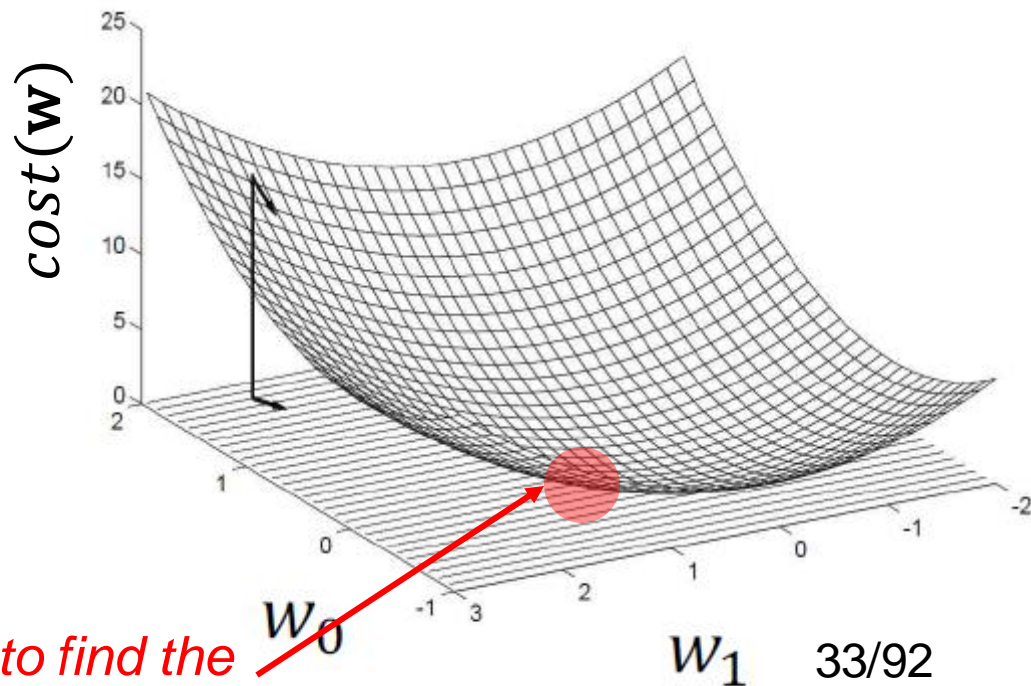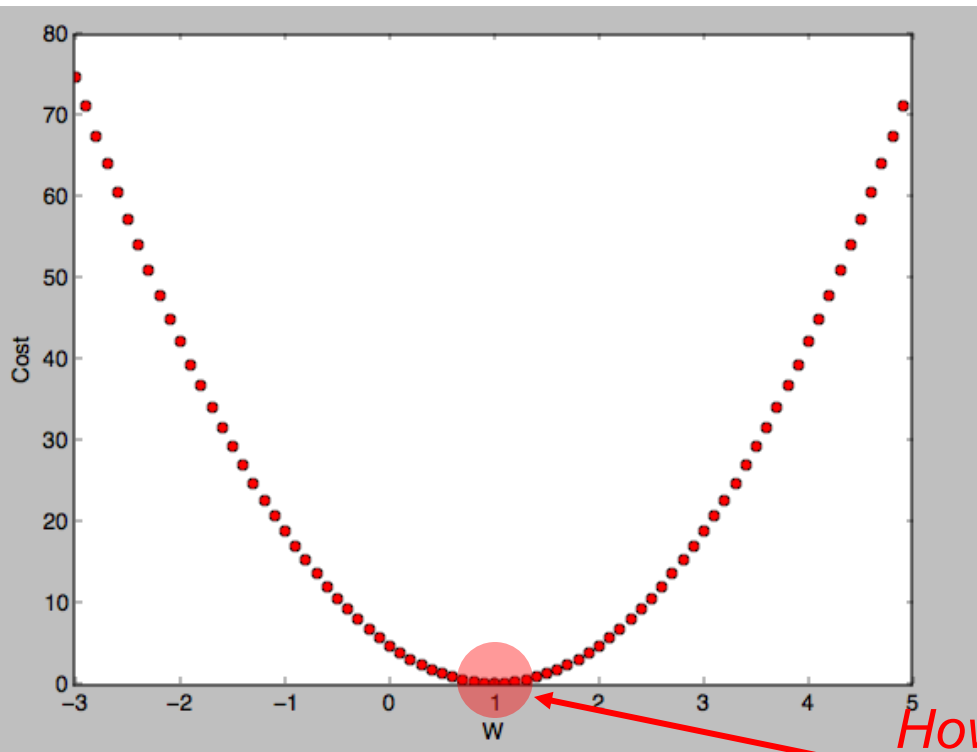  $\frac{1}{2}\left((1 - 1 \times 1)^2 + (2 - 1 \times 2)^2 + (3 - 1 \times 3)^2\right)$

# How cost(W) looks like?

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} (y^{(d)} - \mathbf{w}^T \mathbf{x}^{(d)})^2$$

| x | y |
|---|---|
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |

- $w = 1, E(\mathbf{w}) = 0$

  $\frac{1}{2}\left((1 - 1 \times 1)^2 + (2 - 1 \times 2)^2 + (3 - 1 \times 3)^2\right)$

- $w = 0, E(\mathbf{w}) = ?$

- $w = 2, E(\mathbf{w}) = ?$

# How cost(W) looks like?

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} (y^{(d)} - \mathbf{w}^T \mathbf{x}^{(d)})^2$$

| x | y |
|---|---|
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |

- $w = 1, E(\mathbf{w}) = 0$

  $\frac{1}{2}((1 - 1 \times 1)^2 + (2 - 1 \times 2)^2 + (3 - 1 \times 3)^2)$

- $w = 0, E(\mathbf{w}) = 4.5$

  $\frac{1}{2}((1 - 0 \times 1)^2 + (2 - 0 \times 2)^2 + (3 - 0 \times 3)^2)$

- $w = 2, E(\mathbf{w}) = 4.5$

  $\frac{1}{2}((1 - 2 \times 1)^2 + (2 - 2 \times 2)^2 + (3 - 2 \times 3)^2)$

# How cost(W) looks like?

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} (y^{(d)} - \hat{y}^{(d)})^2$$



Minimum point = 1

# How to minimize cost?

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{d \in D} (y^{(d)} - \hat{y}^{(d)})^2$$



*How to find the minimum point?*

# Gradient Descent Algorithm

- Minimize cost function

- Gradient descent is used for many minimization problems

- For a given cost function, $E(\mathbf{w})$, it will find w
  to minimize cost

- Repeat until you converge to a local minimum

# Gradient Descent Algorithm

## How it works?



1. Start with initial guesses
   - Start at random value

2. Each weight is updated by taking a step into the opposite direction of the gradient $\Delta w_i = -\eta \times \dfrac{\partial E}{\partial w_i}$
   - Compute the partial derivative of the cost function $\dfrac{\partial E}{\partial w_i}$ for each weight

3. Repeat until you converge to a local minimum

# Gradient Descent Algorithm



$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$\eta$ : learning rate (*e.g.* 0.001)

$\frac{\partial E}{\partial w_i}$ : gradient

# Gradient descent algorithm

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d \left(y^{(d)} - \hat{y}^{(d)}\right)^2$$

$$= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} \left(y^{(d)} - \hat{y}^{(d)}\right)^2$$

$$= \frac{1}{2} \sum_d 2\left(y^{(d)} - \hat{y}^{(d)}\right) \frac{\partial}{\partial w_i} \left(y^{(d)} - \hat{y}^{(d)}\right)$$

$$= \sum_d \left(y^{(d)} - \hat{y}^{(d)}\right) \frac{\partial}{\partial w_i} \left(y^{(d)} - w_0 x_0 - w_1 x_1 - \cdots - w_i x_i\right)$$

$$= \sum_d \left(y^{(d)} - \hat{y}^{(d)}\right) \times (-x_i) \qquad \Delta w_i = \eta \times \sum_{d \in D} \left(y^{(d)} - \hat{y}^{(d)}\right) \times (-x_i)$$

# Delta Rule based on Gradient Descent Algorithm (Summary)

$$\hat{y} = w_0 1 + w_1 x_1 + \cdots + w_n x_n$$

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_d \left( y^{(d)} - \hat{y}^{(d)} \right)^2$$

Squared loss function

$$w_i^{new} \leftarrow w_i^{old} + \Delta w_i$$

$$\Delta w_i = -\eta \times \sum_{d \in D} \left( y^{(d)} - \hat{y}^{(d)} \right) \times x_i$$

# Hardware implementations



Frank Rosenblatt, ~1957: Perceptron



Widrow and Hoff, ~1960: Adaline/Madaline

# False Promises

**The New York Times**

## NEW NAVY DEVICE LEARNS BY DOING

*July 8, 1958*

"The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence... Dr. Frank Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers"

Reference: 'NewYork Times' newspaper on July 8, 1958

# (Simple) AND/OR problem: linearly separable?



| x₁ | x₂ | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# Multi-Layer Perceptron (by M. Minsky)



Input Layer

Hidden Layer

Output Layer

$\left(\Sigma\right)$ : Perceptron

2-Layer Neural Network

input layer

hidden layer

output layer

input layer

hidden layer 1    hidden layer 2

output layer

# Multi-Layer Perceptron: Limitation



$$\hat{y} = w_4 \times \overbrace{(w_0 x_0 + w_1 x_1)}^{A} + w_5 \times \overbrace{(w_2 x_1 + w_3 x_1)}^{B}$$

$$= w_4 w_0 x_0 + w_4 w_1 x_1 + w_5 w_2 x_1 + w_5 w_3 x_1$$

$$= (w_4 w_0 + w_5 w_2)x_0 + (w_4 w_1 + w_5 w_3)x_1$$

$$= \underline{(w_4 w_0 + w_5 w_2)} + \underline{(w_4 w_1 + w_5 w_3)}x_1$$

$$\hat{y} = w_a + w_b x_1$$

Still Linear equation
(Line, plane, or hyper-plane)

# Multi-Layer Perceptron: Limitation



$\hat{y} = f(z) = z$

Linear function

# Multi-Layer Perceptron: Activation Function



Non-linear function

Sigmoid function

$$f(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$
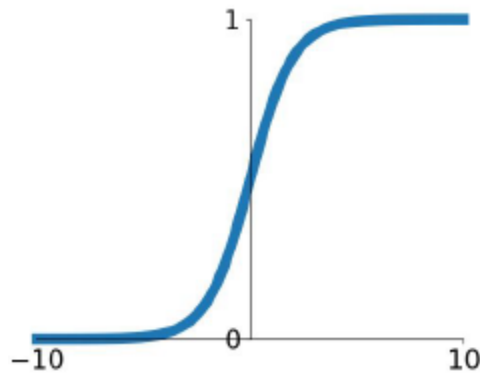
# Activation Functions

Inputs  Weights  Net input function  <span style="color:red">Activation function</span>

**Linear**

$1$  $w_0$

$x_1$  $w_1$

$\vdots$

$x_n$  $w_n$

$\sum x_i w_i$  $z$  $\hat{y}$

$$f(z) = z$$

**Non-linear**

$1$  $w_0$

$x_1$  $w_1$

$\vdots$

$x_n$  $w_n$

$\sum x_i w_i$  $z$  $\hat{y}$

Note: $\frac{d}{dz} f(z) = f(z)(1 - f(z))$  $f(z) = \dfrac{1}{1 + e^{-z}}$

# Activation Functions

**Sigmoid**

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

**Leaky ReLU**

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

**tanh**

$$\tanh(x)$$

**Maxout**

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

**ReLU**
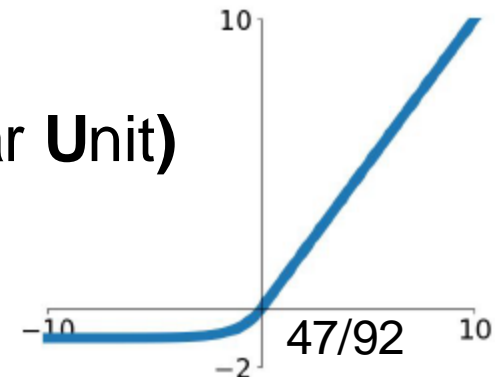**(Re**ctified **L**inear **U**nit**)**
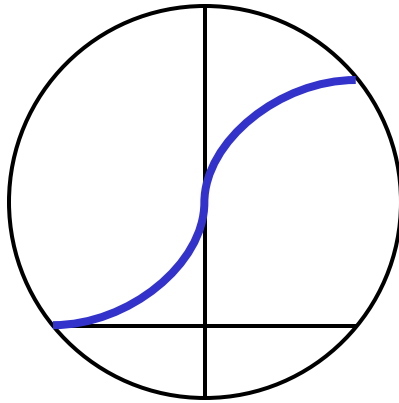
$$\max(0, x)$$

**ELU**
**(E**xponential **L**inear **U**nit**)**

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

# Derivative of Sigmoid Function

$$f(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$
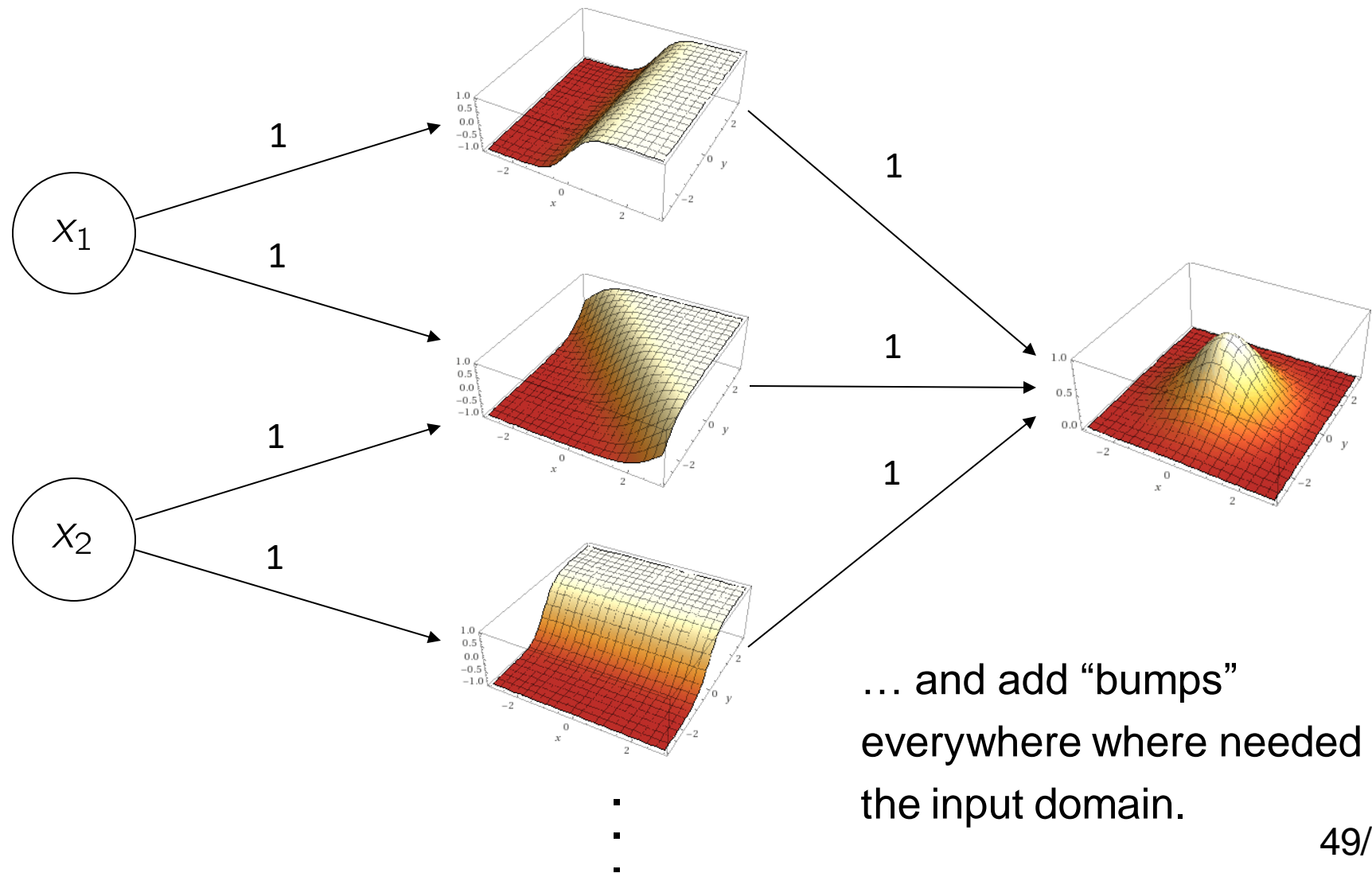
$$\frac{d}{dx}f(x) = f(x)(1 - f(x))$$

$$\frac{d}{dx}f(x) = \frac{d}{dx}\frac{1}{1 + e^{-x}} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$\left(\frac{d}{dx}e^x = e^x\right) \qquad = \frac{1 - 1 + e^{-x}}{(1 + e^{-x})^2} = \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(1 + e^{-x})^2}$$
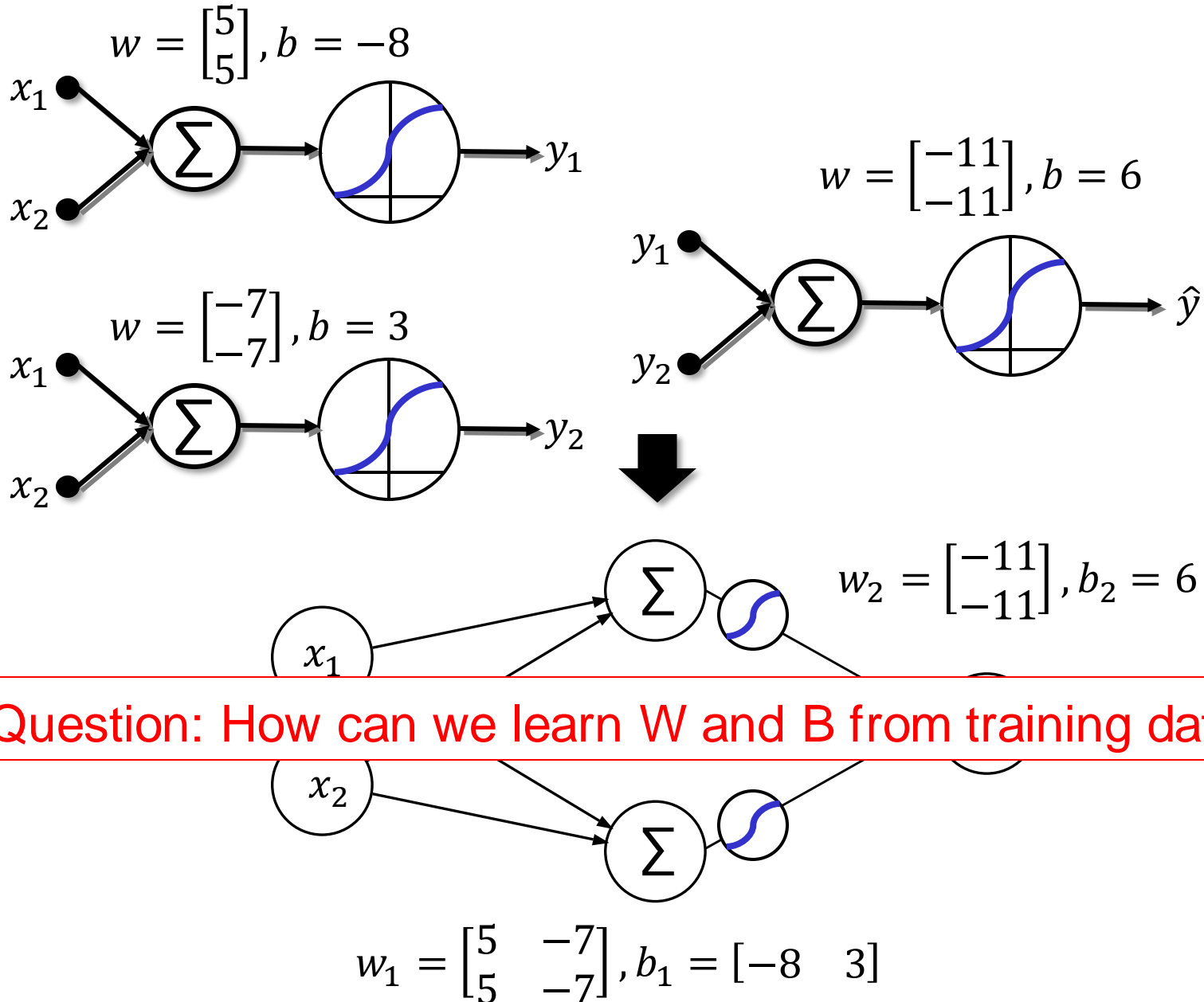
$$= \left(\frac{1}{1 + e^{-x}}\right) - \left(\frac{1}{1 + e^{-x}}\right)^2 = \left(\frac{1}{1 + e^{-x}}\right)\left(1 - \frac{1}{1 + e^{-x}}\right)$$

$$= f(x)(1 - f(x))$$

# Multi-layer Perceptron is Universal Approximator
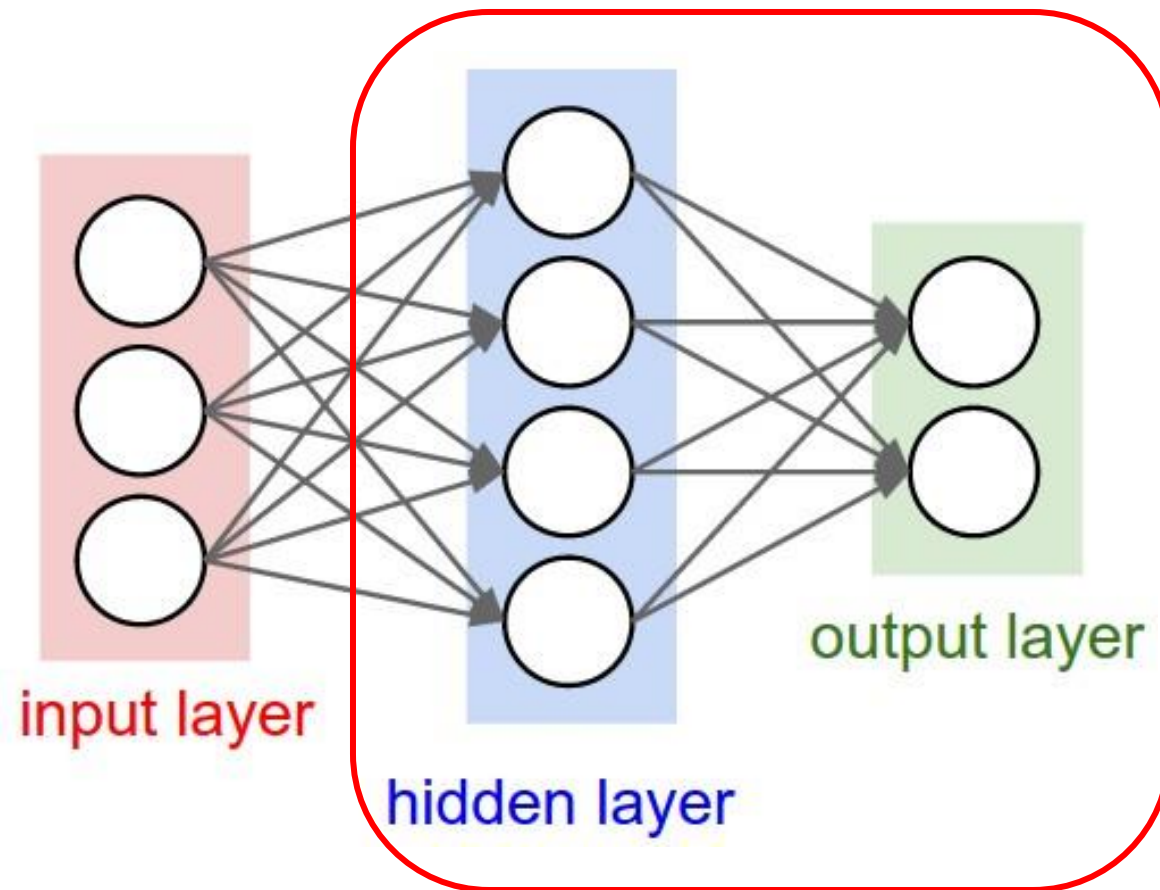
"Proof" by construction:



… and add "bumps" everywhere where needed in the input domain.

# Multi-Layer Perceptron: Forward Propagation



$$w = \begin{bmatrix} 5 \\ 5 \end{bmatrix}, b = -8$$

$x_1$

$x_2$

$y_1$

$$w = \begin{bmatrix} -11 \\ -11 \end{bmatrix}, b = 6$$

$y_1$

$\hat{y}$

$$w = \begin{bmatrix} -7 \\ -7 \end{bmatrix}, b = 3$$

$x_1$

$x_2$

$y_2$

$y_2$

$$w_2 = \begin{bmatrix} -11 \\ -11 \end{bmatrix}, b_2 = 6$$

$x_1$

Question: How can we learn W and B from training data?

$x_2$

$$w_1 = \begin{bmatrix} 5 & -7 \\ 5 & -7 \end{bmatrix}, b_1 = \begin{bmatrix} -8 & 3 \end{bmatrix}$$

# Learning on Multi-Layer Perceptron



$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d \left( y^{(d)} - \hat{y}^{(d)} \right)^2$$

# Learning on Multi-Layer Perceptron



Difficult to calculate $\dfrac{\partial E}{\partial w_j}$
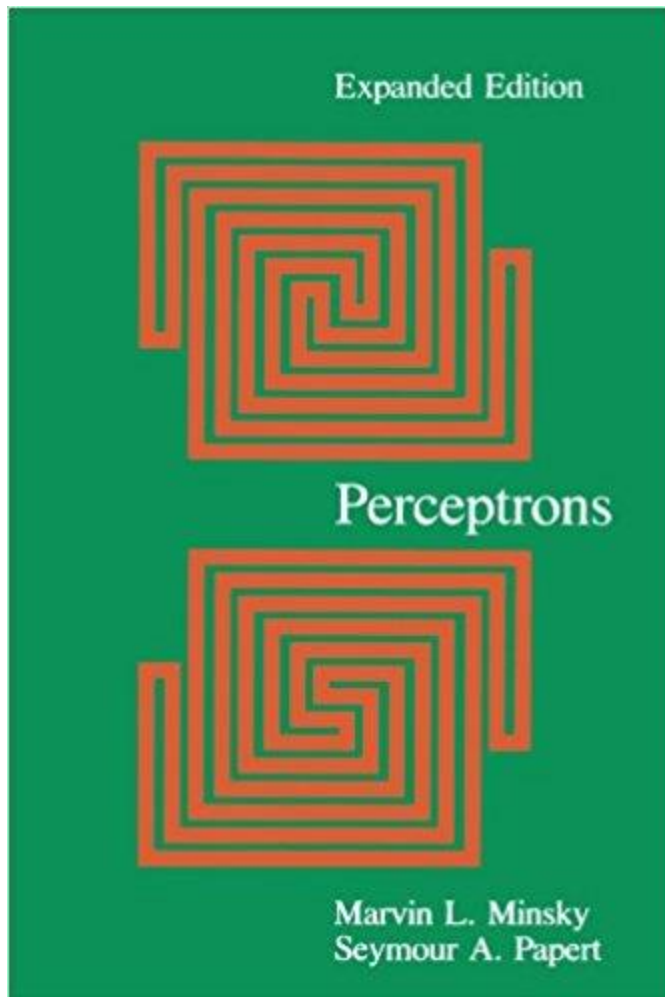
# Learning on Multi-Layer Perceptron



Unsolved problem for 20 years

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d \left( ? - \hat{y}^{(d)} \right)^2$$

# Limitation of Multi-layer Perceptron
## By Marvin Minsky, founder of the MIT AI Lab.



- We need to use MLP, multilayer perceptrons (multilayer neural nets)

- No one on earth had found a viable way to train MLPs good enough to learn such simple functions

# Backpropagation
# (1974, 1982 by Paul Werbos, 1986 by Hinton)

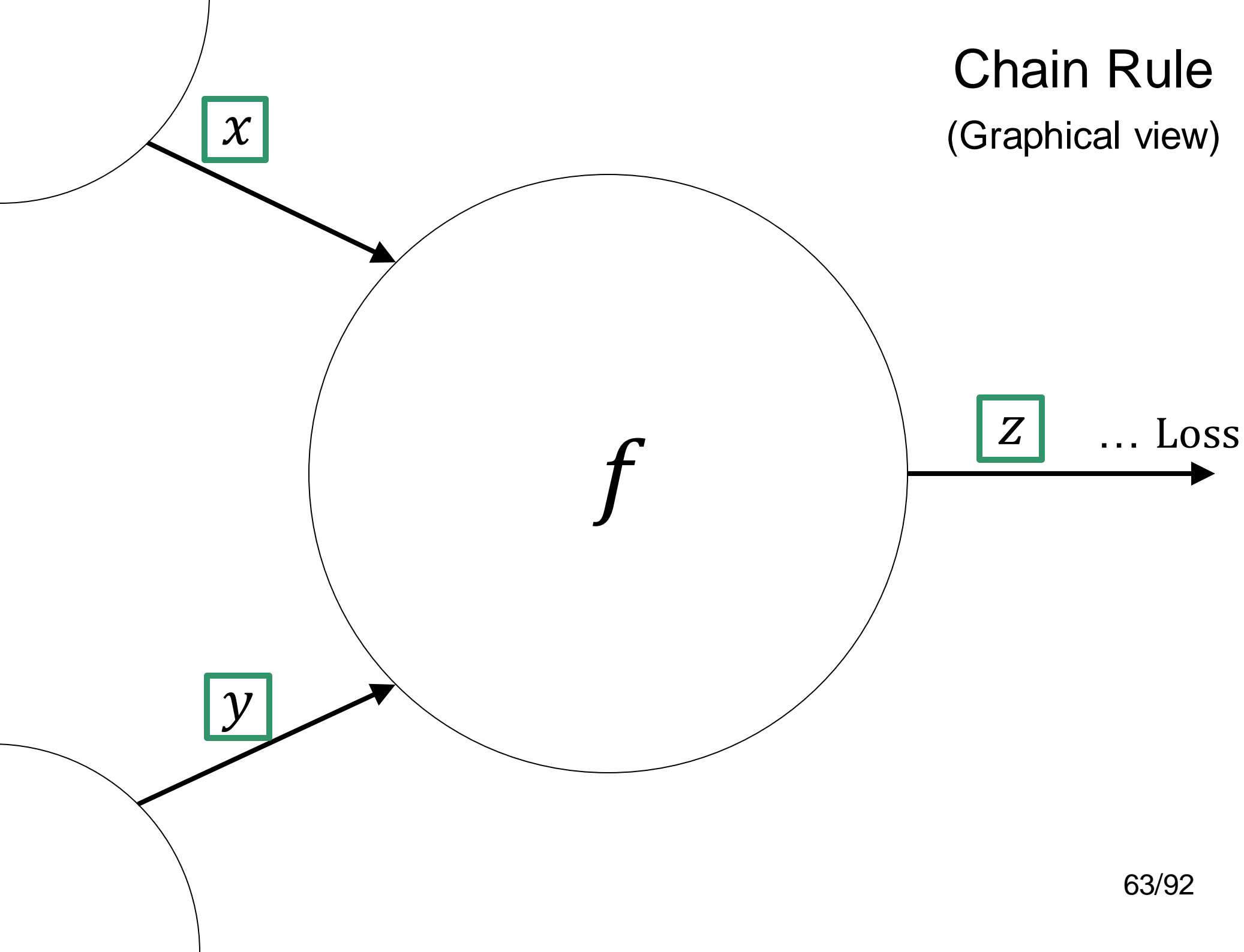https://devblogs.nvidia.com/inference-next-step-gpu-accelerated-deep-learning/

# Chain Rule

$$f = \underline{f}(g); \; g = \underline{g}(x)$$

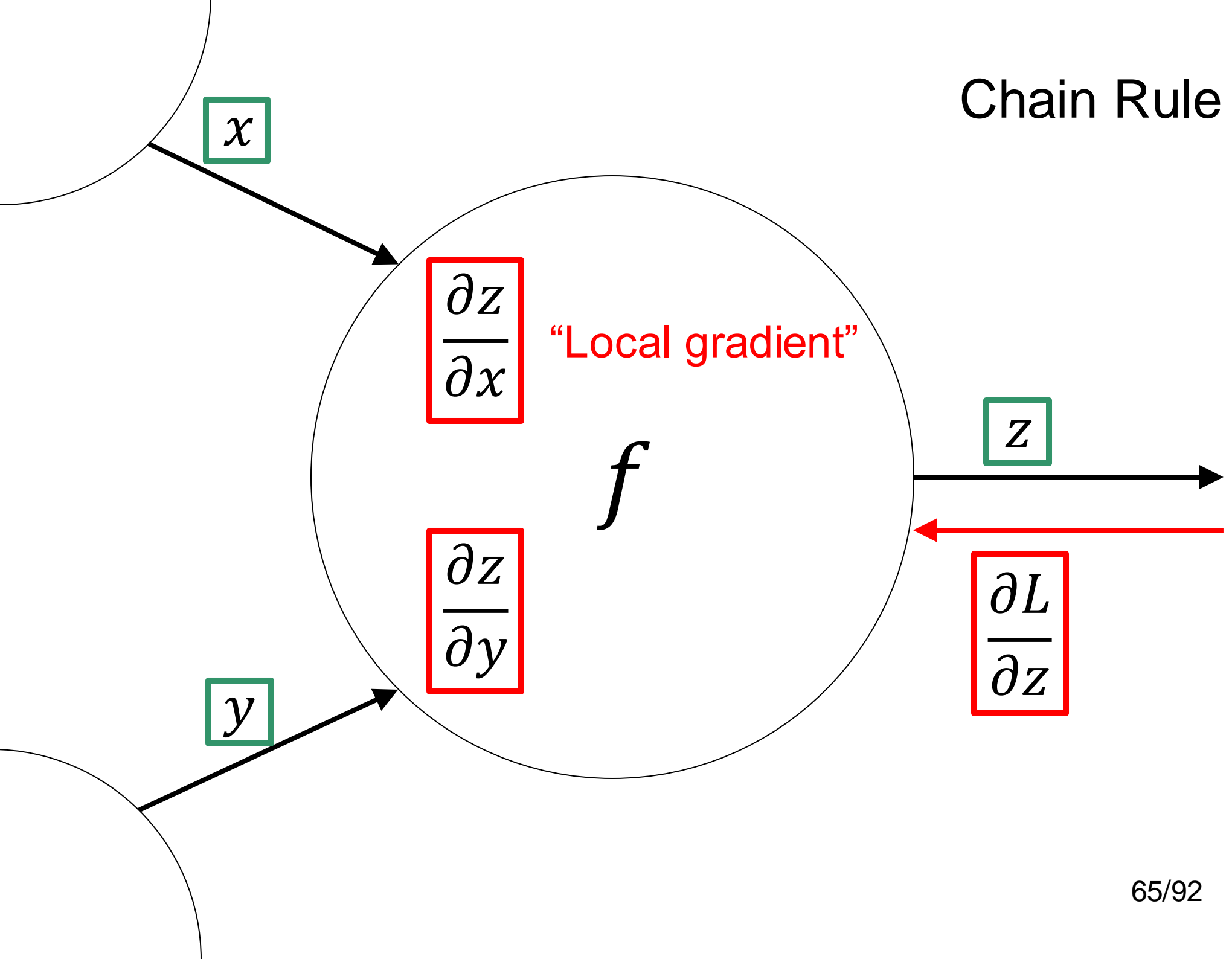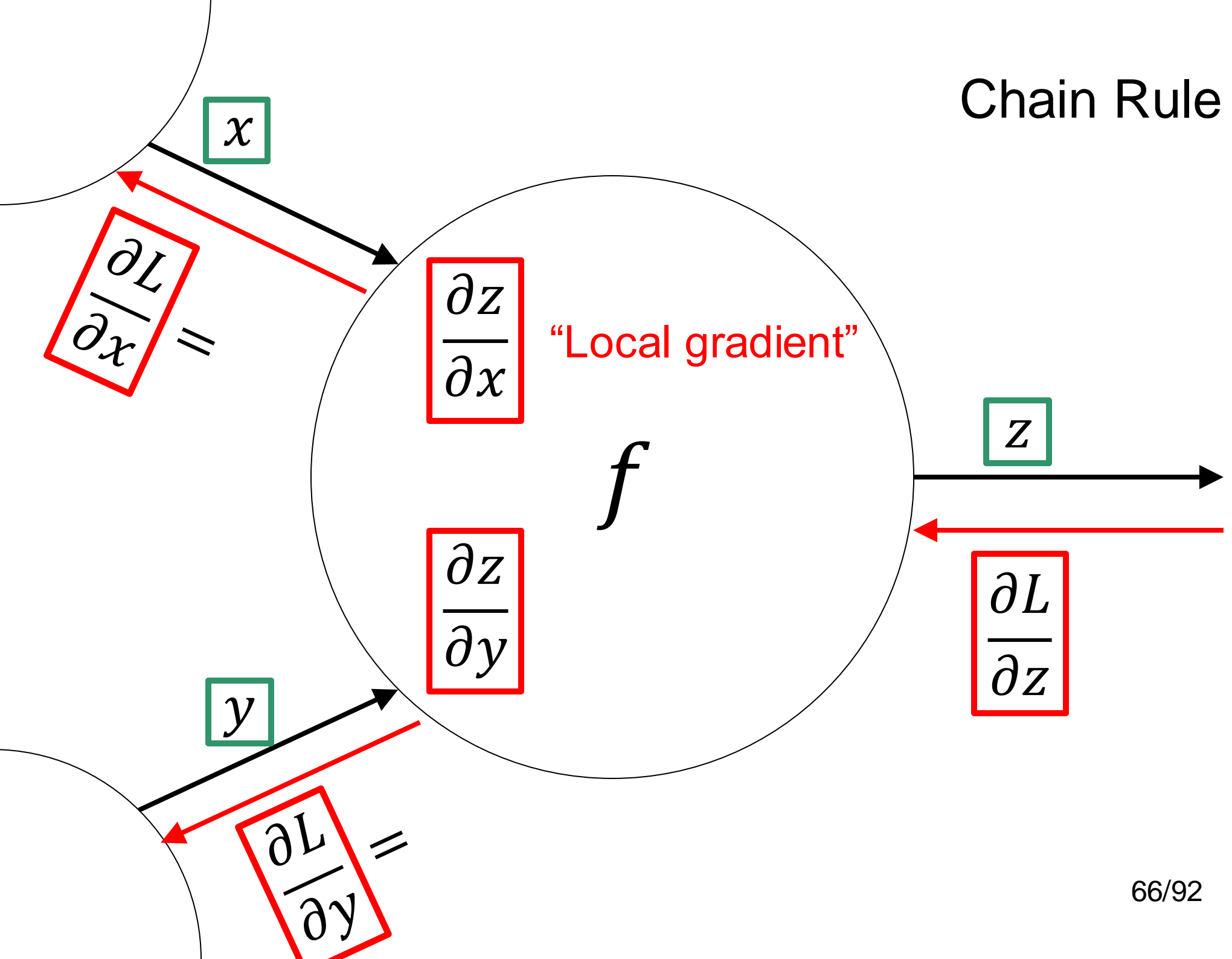$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

https://mathinsight.org/image/functino_machines_composed

# Chain Rule
(Graphical view)

$x$

$f$

$z$  ... Loss

$y$

Chain Rule

$x$

$\dfrac{\partial z}{\partial x}$ "Local gradient"

$f$

$\dfrac{\partial z}{\partial y}$

$y$

$z$

$x$

$y$

$f \times$

$z$ ... Loss

$x$ = 2

$\dfrac{\partial z}{\partial x}$   "Local gradient"

$f$ ×

$z$ = 6

$\dfrac{\partial z}{\partial y}$

$y$ = 3

$x$ = 2

"Local gradient"

$$\frac{\partial z}{\partial x} = \frac{\partial x \cdot y}{\partial x} = y$$

$$f \times$$

$$\frac{\partial z}{\partial y} = \frac{\partial x \cdot y}{\partial x} = x$$

$z$ = 6

$y$ = 3

$x$ = 2

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$
$$= 5 \cdot y = 15$$

$\frac{\partial z}{\partial x}$

"Local gradient"

$f \times$

$z$ = 6

$\frac{\partial z}{\partial y}$

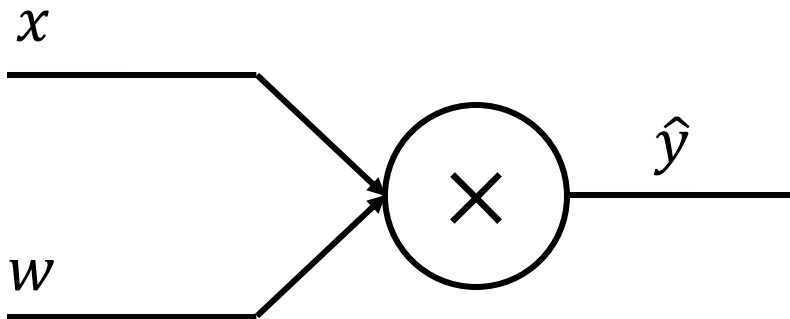$\frac{\partial L}{\partial z}$ = **5**

$y$ = 3

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$
$$= 5 \cdot x = 10$$

# Computational graph

$$\hat{y} = x \times w$$
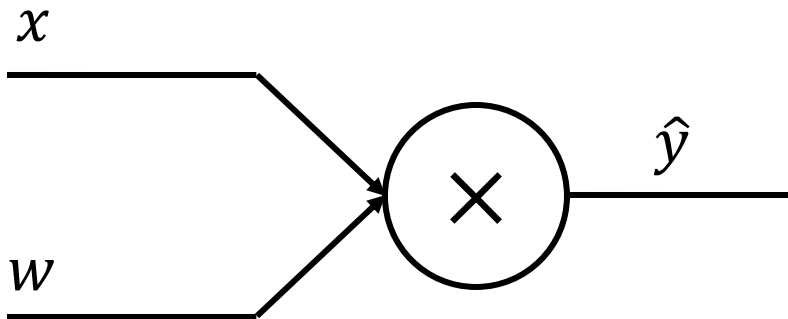
# Computational graph

$$\hat{y} = x \times w$$

# Computational graph

$$\hat{y} = x \times w \qquad loss = (\hat{y} - y)^2 = (x \times w - y)^2$$

$$\frac{1}{2}\sum_{d \in D}\left(\hat{y}^{(d)} - y^{(d)}\right)^2$$

**Our loss function**

# Computational graph

$$\hat{y} = x \times w \qquad loss = (\hat{y} - y)^2 = (x \times w - y)^2$$

$$\hat{y} = x \times w \qquad loss = (\hat{y} - y)^2 = (x \times w - y)^2$$



$x$ = 1

$\hat{y}$

$s$

$loss$ = ?

$\times$

$-$

$^2$

$w$ = 1

$y$ = 2

$$\hat{y} = x \times w \qquad loss = (\hat{y} - y)^2 = (x \times w - y)^2$$

$x$ = 1

$w$ = 1

$\hat{y}$ = 1

× 

$S$ = -1

−

^2

$loss$ = ?

$y$ = 2

$$\hat{y} = x \times w \quad loss = (\hat{y} - y)^2 = (x \times w - y)^2$$



$x = 1$

$w = 1$

$\hat{y} = 1$

$S = -1$

$^\wedge 2$

$loss = 1$

$y = 2$

$$\hat{y} = x \times w \qquad loss = (\hat{y} - y)^2 = (x \times w - y)^2$$

$$\frac{\partial loss}{\partial w} = \frac{\partial loss}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w} = -2 \cdot x = -2 \cdot 1 = -2$$

$$\frac{\partial xw}{\partial w} = x$$

$x$ = 1

$w$ = 1

×

$\hat{y}$ = 1

−

$s$ = -1

^2

$loss$ = 1

$$\frac{\partial \hat{y} - y}{\partial \hat{y}} = 1$$

$y$ = 2

$$\frac{\partial s^2}{\partial s} = 2s$$

$$\frac{\partial loss}{\partial w} = -2$$

$$\frac{\partial loss}{\partial \hat{y}} = \frac{\partial loss}{\partial s} \frac{\partial s}{\partial \hat{y}} = -2 \cdot 1 = -2$$

$$\frac{\partial loss}{\partial s} = 2s = -2$$

# Backpropagation: Exercise



$x$ = 2

$w$ = 1

$\hat{y}$

$s$

$loss$

×

−

^2

$y$ = 4

$$\frac{\partial loss}{\partial w} = ?$$

# Backpropagation: Summary

1:  Initialize all weights to small random values.
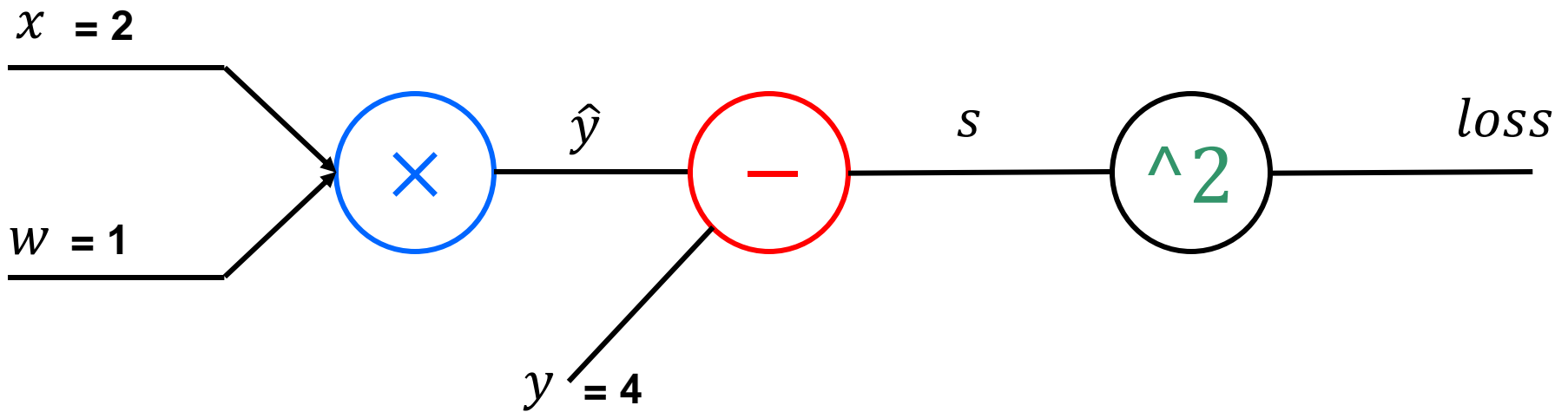
2:  **repeat**

3:      **for** each training example do

4:          <u>Forward propagate</u> the input features of the example

        to determine the MLP's outputs.

5:          <u>Back propagate</u> <u>error</u> to generate $\Delta w_i$ for all weights $\Delta w_i$

6:          <u>Update</u> the weights using $\Delta w_i$

7:      **end for**

8:  **until** stopping criteria reached.

# Summary

- We learned what a perceptron and multilayer perceptron is

- We have some intuition about using gradient descent on an error function

- We know a learning delta rule for updating weights in order to minimize the error: $\Delta w_i = -\eta \times \frac{\partial E}{\partial w_i}$

- We know activation function for non-linearity

- We can use this rule to learn an MLP using the backpropagation algorithm

# Optimizer: How to Optimize?

## Gradient Descent (GD) vs. Stochastic GD (SGD)

<u>Objective to Minimize:</u>

$$E(\mathbf{w}) \equiv \frac{1}{D} \sum_{d \in D} E^{(d)}(\mathbf{w})$$

<u>GD:</u>

**while** True:

$$\mathbf{w}^{new} \leftarrow \mathbf{w}^{old} - \eta \frac{\partial}{\partial \mathbf{w}} \frac{1}{D} \sum_{d \in D} E^{(d)}(\mathbf{w})$$

$$O(N)$$

<u>SGD:</u>

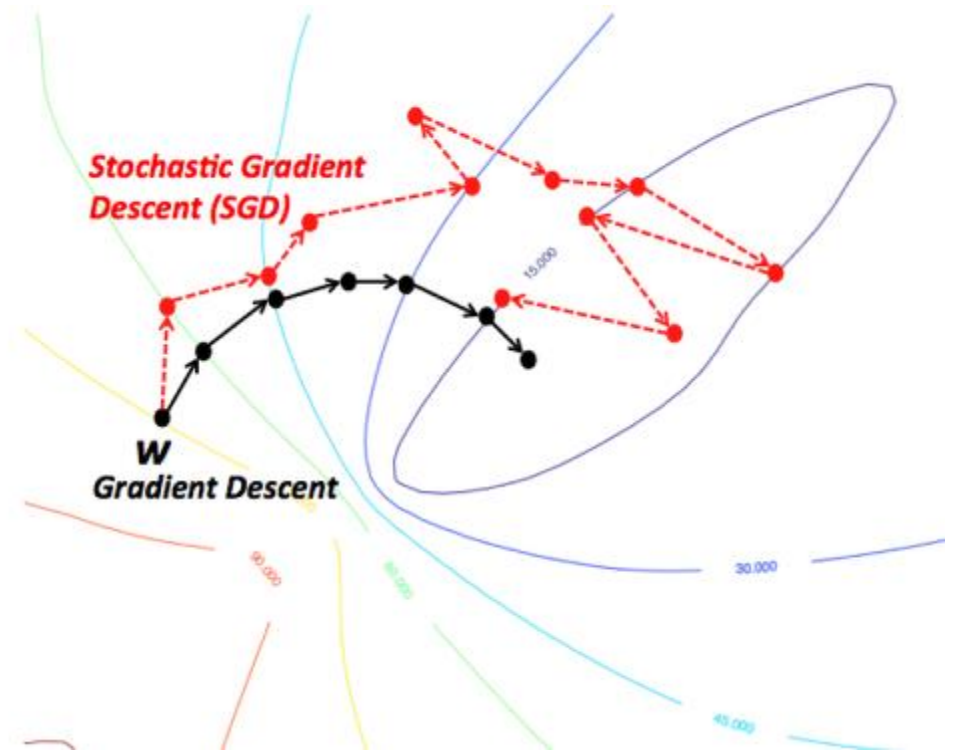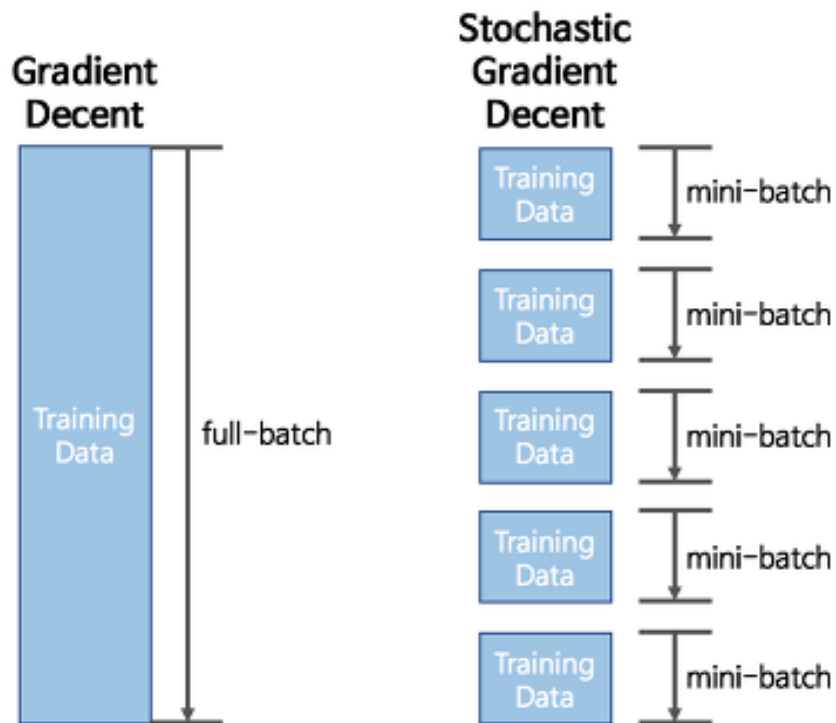**while** True:

$$d \leftarrow mini\ batch(1,N)$$

$$\mathbf{w}^{new} \leftarrow \mathbf{w}^{old} - \eta \frac{\partial E^{(d)}(\mathbf{w})}{\partial \mathbf{w}}$$
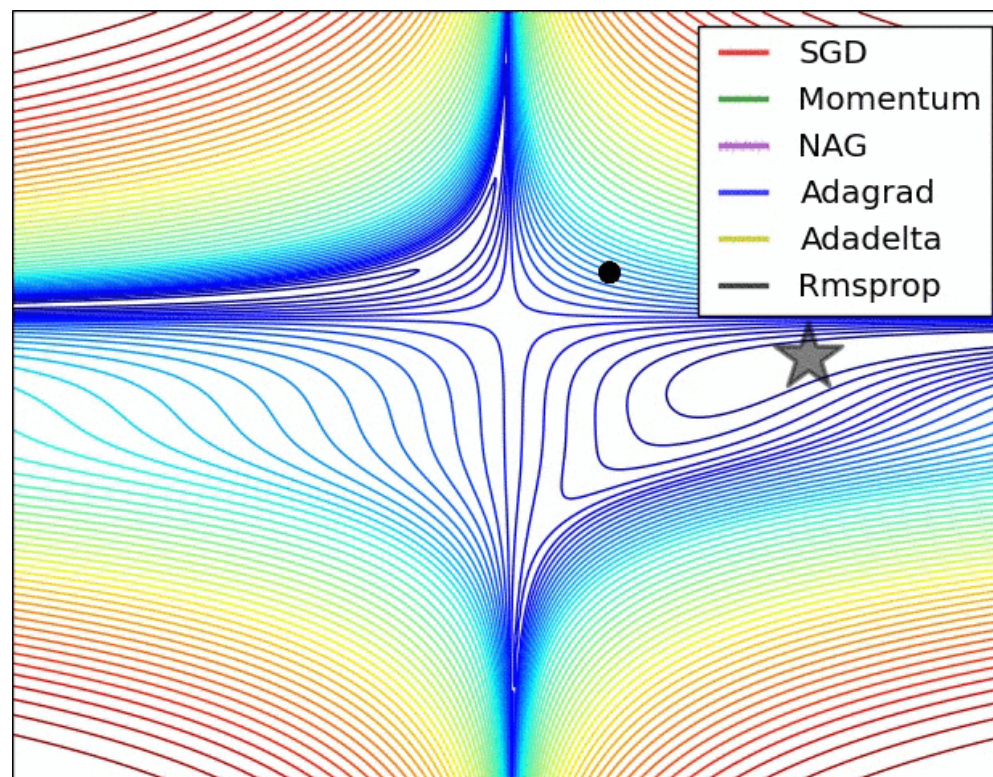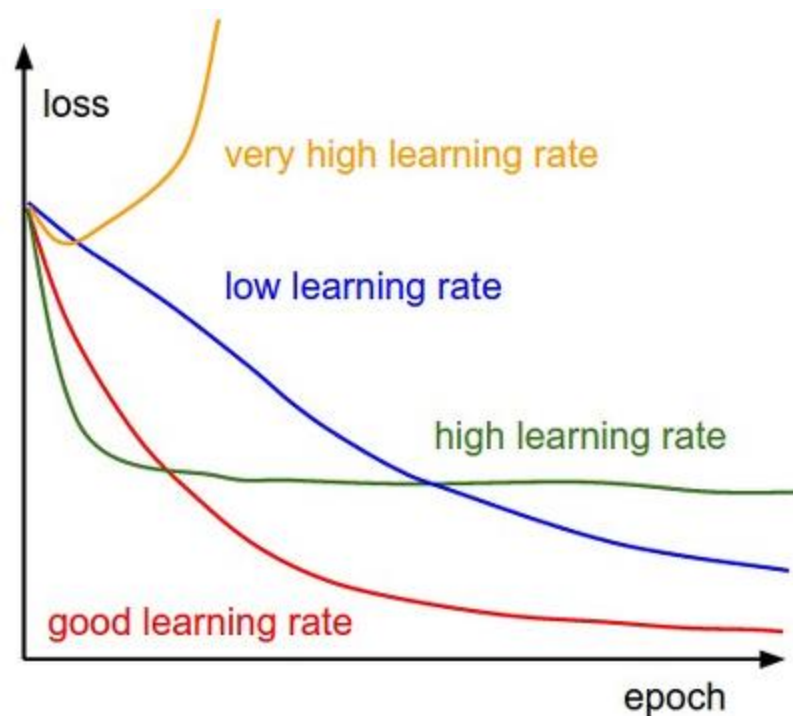
$$O(d)$$

# Optimizer: How to Optimize?
## Gradient Descent vs. SGD



GD vs. SGD

# Limitations of SGD and its Solutions

http://ruder.io/optimizing-gradient-descent/

# Tutorial on Neural Network: Zero to All !

- Perceptron from scratch

  - Based and basic algebra

- Perceptron based on pytorch

- Multi-layer perceptron

**PYTORCH**

| Step | Method | Data | Model | Forward | Loss | Backward | Update |
|------|--------|------|-------|---------|------|----------|--------|
| 1 | Perceptron | x | x | x | x | x | x |
| 2 | | x | x | x | x | o | x |
| 3 | | x | o | o | o | o | o |
| 4 | MLP | o | o | o | o | o | o |

# Deep Neural Network Toolbox



Artificial Intelligence and Machine Learning SkillsFuture Courses and Training

- Why **PYTORCH**?

  - More pythonic

    - Flexible, intuitive and cleaner code, and easy to debug

- More neural networkic

  - Write code as the network works and forward/backward

# Multi-Layer Perceptron: Programming
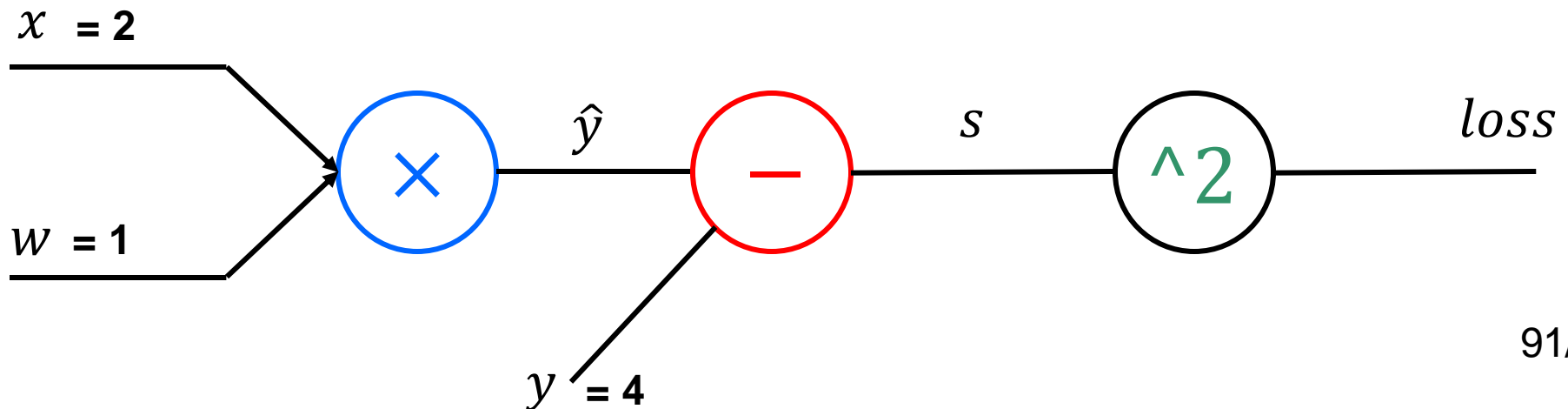
## 1. Training data

| x | y |
|---|---|
| 1 | 2 |
| 2 | 4 |
| 3 | 6 |

## 2. Model (forward pass)

$$\hat{y} = x \times w$$

## 3. Loss function

$$loss = (\hat{y} - y)^2 = (x \times w - y)^2$$

$x = 2$

$w = 1$

$\hat{y}$

$s$

$\wedge 2$

$loss$

$y = 4$

# Multi-Layer Perceptron: Programming

## 1. Training data

| x | y |
|---|---|
| 1 | 2 |
| 2 | 4 |
| 3 | 6 |

## 2. Model (forward pass)

$$\hat{y} = x \times w$$

## 3. Loss function

$$loss = (\hat{y} - y)^2 = (x \times w - y)^2$$

**Learning Process (Very important !!!!)**

Forward → Loss → Backward → Update →