*Statistical Machine Learning Theory*

# Neural Networks

Hisashi Kashima

kashima@i.Kyoto-u.ac.jp
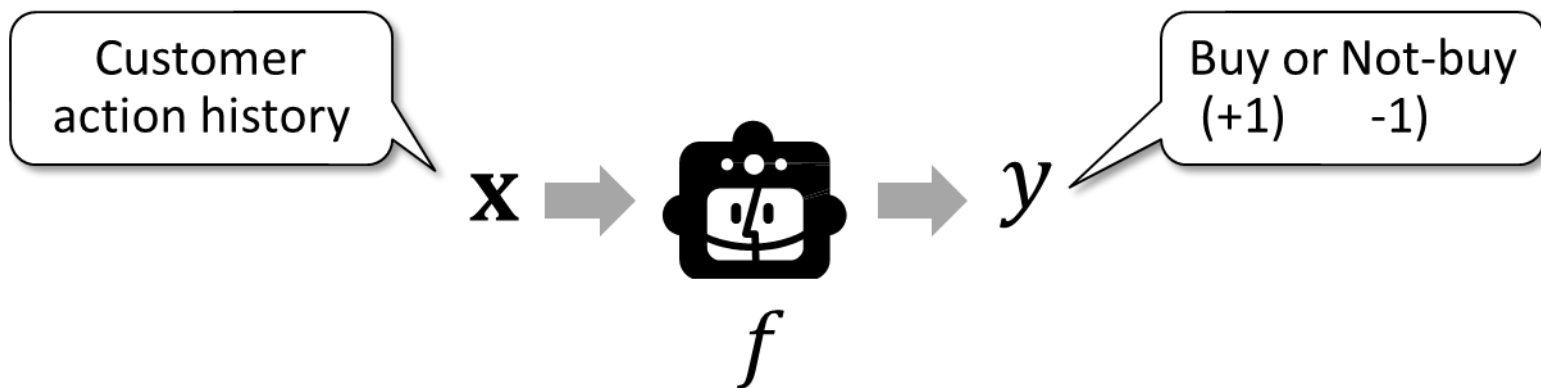
# Logistic Regression

# Classification:
## Supervised learning for predicting discrete variable

- Goal: Obtain a function $f: \mathcal{X} \rightarrow \mathcal{Y}$

  - Input domain: $\mathcal{X} = \mathbb{R}^D$

  - $\mathcal{Y}$: discrete domain

    - We focus on two-class classification: $\mathcal{Y} = \{+1, -1\}$

- Training dataset: $N$ pairs of an input and an output
  $$\{(\mathbf{x}^{(1)}, y^{(1)}), \ldots, (\mathbf{x}^{(N)}, y^{(N)})\}$$

Customer action history

$\mathbf{x} \rightarrow f \rightarrow y$
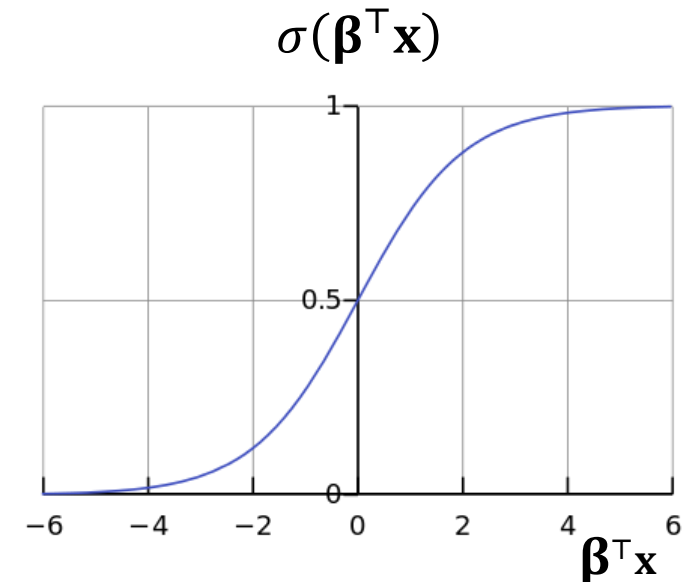
Buy or Not-buy
(+1)     -1)

# Logistic regression:
# A probabilistic model for binary classification

- Logistic regression model give the conditional probability

$$f_{\mathbf{w}}(y = +1|\mathbf{x}) = \sigma(\mathbf{w}^{\top}\mathbf{x}) = \frac{1}{1+\exp(-\mathbf{w}^{\top}\mathbf{x})}$$

- Logistic (sigmoid) function $\sigma: \mathbb{R} \rightarrow (0,1)$

  - That converts real numbers to "probabilities"

$\sigma(\boldsymbol{\beta}^{\top}\mathbf{x})$



$\boldsymbol{\beta}^{\top}\mathbf{x}$

## Objective function to train a logistic regression model

- Cross entropy to minimize:

For positive class data

$$L(\mathbf{w}) = \sum_{i=1}^{N} \delta(y^{(i)} = 1)\log f_{\mathbf{w}}(y^{(i)} = +1|\mathbf{x}^{(i)})$$

$$+ \sum_{i=1}^{N} \delta(y^{(i)} = -1)\log(1 - f_{\mathbf{w}}(y^{(i)} = +1|\mathbf{x}^{(i)}))$$

$$= \sum_{i=1}^{N} \log(1 + \exp(-y^{(i)}\mathbf{w}^{\top}\mathbf{x}^{(i)}))$$

For negative class data

- Equivalent to
  - Logistic loss: upper bound of 0-1 loss (#mistakes)
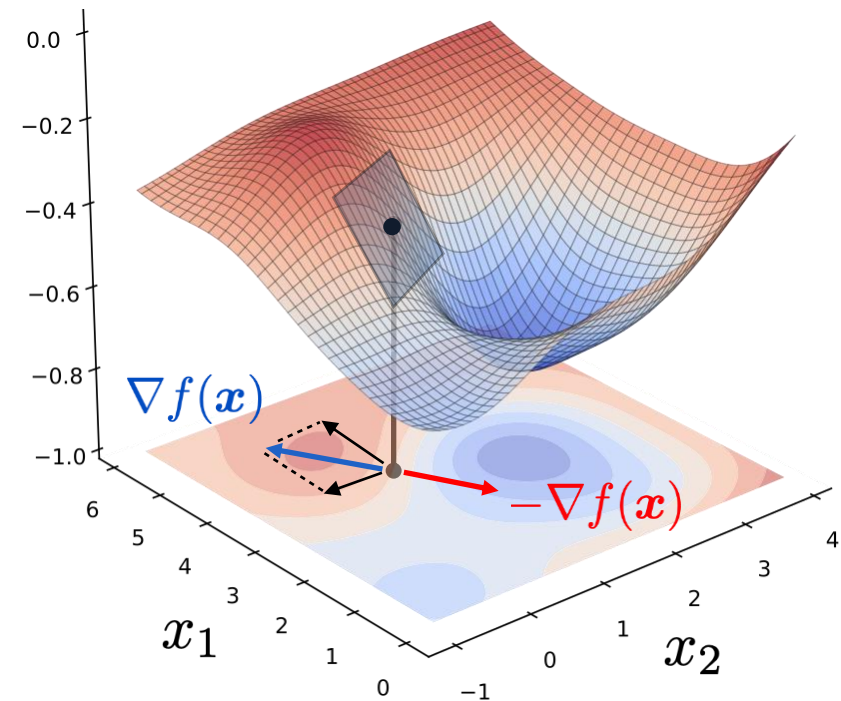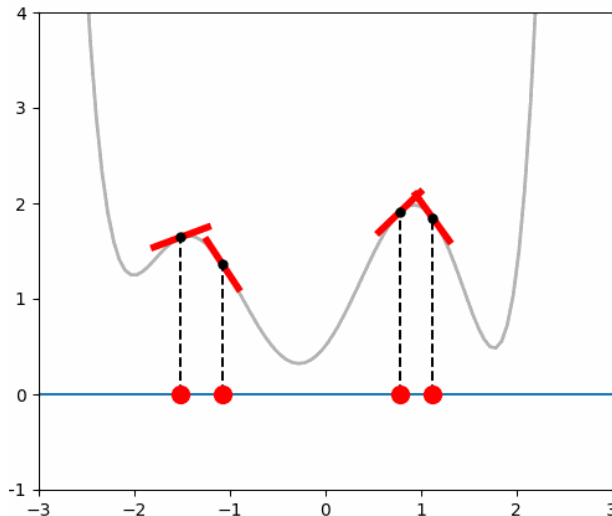  - Negative log-likelihood (for maximum likelihood estimation)

# Gradient descent:
## A simplest parameter optimization method

- Iteratively refine the current parameter $\mathbf{w}$ to $\mathbf{w}^{\text{NEW}}$:

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} - \eta \nabla L(\mathbf{w})$$

- Gradient $\nabla L(\mathbf{w})$ is the most steepest direction of the objective function $L(\mathbf{w})$

- "learning rate" $\eta$

# Gradient descent for logistic regression: Gradient of objective function is all you need
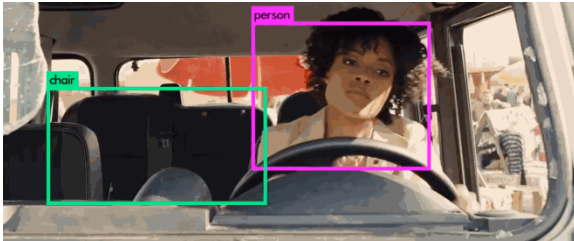
- Once we obtain the gradient, we can apply GD

  - Obj. func.: $L(\mathbf{w}) = \sum_{i=1}^{n} \ln\left(1 + \exp\left(-y^{(i)} \mathbf{w}^{\top} \mathbf{x}^{(i)}\right)\right)$

  - Gradient: $\dfrac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^{n} \dfrac{y^{(i)} \mathbf{x}^{(i)}}{1 + \exp\left(\mathbf{w}^{\top} \mathbf{x}^{(i)}\right)}$

  - GD update: $\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} - \eta \sum_{i=1}^{n} \dfrac{y^{(i)} \mathbf{x}^{(i)}}{1 + \exp\left(\mathbf{w}^{\top} \mathbf{x}^{(i)}\right)}$

- Approximation using only one data instance

  - Stochastic gradient descent (SGD):

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} - \eta \frac{y^{(i)} \mathbf{x}^{(i)}}{1 + \exp\left(\mathbf{w}^{\top} \mathbf{x}^{(i)}\right)}$$

KYOTO UNIVERSITY

# Neural Networks

# Success of "deep learning" :
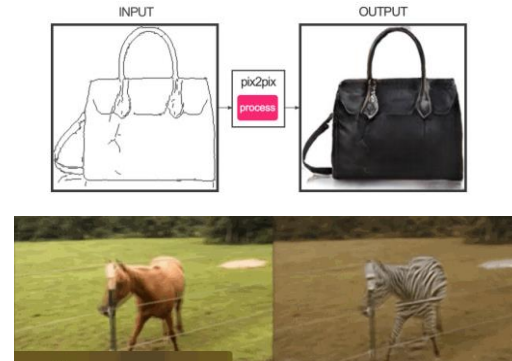## Real world applications

Image recognition

Machine Translation

Image/movie transformation

"Deep Fake"



AlphaGo

AlphaFold2
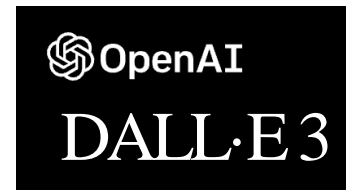
AlphaTensor

AlphaGo

AlphaFold

AlphaTensor

"Generative" AI

# Neural networks:
## Multi-layered logistic regression models

- (Very roughly speaking, ) a neural network is multi-layered logistic regression models

  - Outputs of some logistic regression models are inputs to other logistic regression models

- The form of final output is still $\Pr(y = +1|\mathbf{x})$

Logistic regression

Neural network (2 layers)

$x_1$ — $w_1$

Logistic function

$x_2$ — $w_2$

$\rightarrow \Pr(y = +1|\mathbf{x})$

$x_1$

$w_{11}$

$w_{12}$

$w_{21}$

$w_{22}$

$w_1$

$w_2$

$x_2$

$\rightarrow \Pr(y = +1|\mathbf{x})$

Layer 1

Layer 2

# Why do we need to stack layers?
# Nonlinear classification

- (1-layer) logistic regression only allow linear classification (AND/OR)

- Gains *non-linear* expressive power by stacking two layers (XOR)

- Universal approximation theorem: neural network can represent arbitrary functions by introducing many intermediate units

Logistic regression

$x_1 \longrightarrow \boxed{\times 2}$

$\boxed{} \rightarrow \Pr(y = +1|\mathbf{x})$

$x_2 \longrightarrow \boxed{\times 1}$

$\Pr(y = +1|\mathbf{x})$

Neural network (2 layers)

$\boxed{\times 2}$
$x_1$
$\boxed{\times 1}$
$x_2$
$\boxed{\times -2}$
$\boxed{\times 1}$

$\boxed{} \rightarrow \boxed{\times 1}$
$\boxed{} \rightarrow \boxed{\times 1}$

$\boxed{} \rightarrow \Pr(y = +1|\mathbf{x})$

$\Pr(y = +1|\mathbf{x})$

# Parameter estimation for neural networks: Gradient is all you need

- The objective function for training a neural network is the same as that of logistic regression, i.e., the *cross entropy*

  - When SGD is used, the cross entropy (only for the $i$-th instance) is

  $$L^{(i)}(\mathbf{w}) = \delta(y^{(i)} = 1)\log f_{\mathbf{w}}(y^{(i)} = +1|\mathbf{x}^{(i)}) + \delta(y^{(i)} = -1)\log f_{\mathbf{w}}(y^{(i)} = -1|\mathbf{x}^{(i)})$$

- Gradient descent update: $\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} - \eta\nabla L^{(i)}(\mathbf{w})$

- How can we obtain the gradient $\nabla L^{(i)}(\mathbf{w}) = \partial L^{(i)}/\partial \mathbf{w}$?

  - E.g., how can we obtain $\partial L^{(i)}/\partial w_{12}$? (They are separated via $w_1$)

# Error back propagation:
## An efficient strategy to compute the gradient

- Once we differentiate the objective function, we can apply SGD

- *Error backpropagation* can compute the gradient

  - Forward path: computes output (objective func.) from the input

  - Backward path: computes derivatives from output to input

    - All derivatives are computed in a recursive manner

# How error back propagation works:
## Efficient computation using the chain rule of derivatives

- 1-dimensional case (of no practical use)



$$x^{(i)} \rightarrow \boxed{w_1} \rightarrow \boxed{\sigma} \rightarrow \boxed{w_2} \rightarrow \boxed{\sigma} \rightarrow f \rightarrow L^{(i)}$$

$$z_1 = w_1 x^{(i)} \qquad z_2 = \sigma(z_1) \qquad z_3 = w_2 z_2 \qquad f = \sigma(z_3)$$

$$L^{(i)}(w_1, w_2) = \delta(y^{(i)} = 1)\log f(x^{(i)}) + \delta(y^{(i)} = -1)\log\left(1 - f(x^{(i)})\right)$$

- Chain rule of derivatives:

$$\bullet \quad \frac{\partial L^{(i)}}{\partial w_2} = \frac{\partial L^{(i)}}{\partial f} \cdot \frac{\partial f}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_2}$$

$$\bullet \quad \frac{\partial L^{(i)}}{\partial w_1} = \frac{\partial L^{(i)}}{\partial f} \cdot \frac{\partial f}{\partial z_3} \cdot \frac{\partial z_3}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$
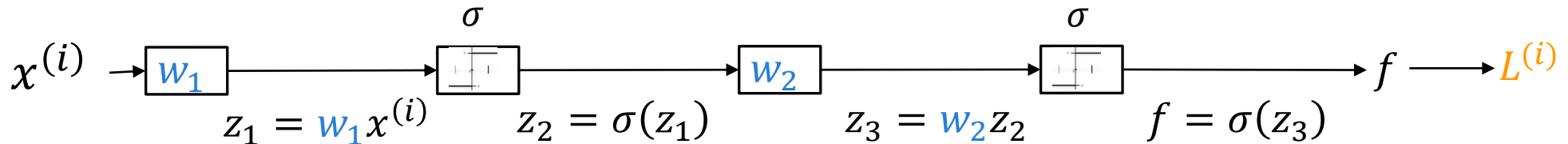
$$\frac{\partial L^{(i)}}{\partial f} = \frac{\delta(y^{(i)} = +1)}{f(x^{(i)})} - \frac{\delta(y^{(i)} = -1)}{1 - f(x^{(i)})}$$

$$\frac{\partial f}{\partial z_3} = \sigma(z_3)(1 - \sigma(z_3))$$

$$\frac{\partial z_3}{\partial w_2} = z_2$$

$$\frac{\partial z_3}{\partial z_2} = w_2$$

$$\frac{\partial z_2}{\partial z_1} = \sigma(z_1)(1 - \sigma(z_1))$$

$$\frac{\partial z_1}{\partial w_1} = x^{(i)}$$

Already obtained in forward path

# How error back propagation works:
## Efficient computation using the chain rule of derivatives

- Naïve application of the chain rules requires $O(|U|^2)$ computation

$$x^{(i)} \rightarrow \boxed{w_1} \xrightarrow{} \boxed{\sigma} \xrightarrow{} \boxed{w_2} \xrightarrow{} \boxed{\sigma} \xrightarrow{c} f \rightarrow L^{(i)}$$

$$z_1 = w_1 x^{(i)} \qquad z_2 = \sigma(z_1) \qquad z_3 = w_2 z_2 \qquad f = \sigma(z_3)$$

- Backward computation allows reuse of common parts
  $\Rightarrow$ results in $O(|U|^2)$ computation     $|U|$: number of operation units

- Chain rule of derivatives:

$$\bullet \quad \frac{\partial L^{(i)}}{\partial w_2} = \frac{\partial L^{(i)}}{\partial f} \cdot \frac{\partial f}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_2}$$

$$\bullet \quad \frac{\partial L^{(i)}}{\partial w_1} = \frac{\partial L^{(i)}}{\partial f} \cdot \frac{\partial f}{\partial z_3} \cdot \frac{\partial z_3}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$
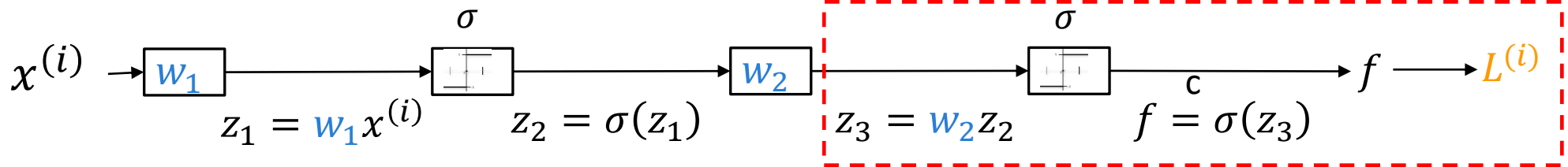
Reusable common part in backward path

# Computational Graphs and Automatic Differentiation

# Computational graph:
## General representation of NN structure as a DAG

- A neural net can be constructed by stacking logistic regression models

    - Parameter estimation requires differentiating the output (the objective function) by parameters in the middle of the NN

    - This can be computed by applying the chain rule on this graph

    $\Rightarrow$ Let us generalize them!

- Computational graph:
  A directed acyclic graph representing computational process from input to output (NN decision process) using simple computational units:

    - Addition and multiplication (matrix multiplication)

    - Sigmoid transformation (simple nonlinear transformation)

    $\Rightarrow$ Can we also use other types of units?

# Requirements for computational unit in NN:
## Output and its derivatives w.r.t. inputs and parameters

- Q: What computational units are allowed to be used in NN?

  A: Required to have its output is differentiable w.r.t its

  1. Inputs
  2. Parameters

- In other words, we can use arbitrary units that offer

  1. Output (for its input) : $\mathbf{y} = f(\mathbf{x}; \boldsymbol{W})$ $\Bigg\}$ Used in forward path

  2. Derivative of the output w.r.t its inputs: $\partial\mathbf{y}/\partial\mathbf{x}$

  3. Derivative of the output w.r.t its parameters : $\partial\mathbf{y}/\partial\boldsymbol{W}$ $\Bigg\}$ Used in backward path

$$\mathbf{x} \longrightarrow \boxed{f} \longrightarrow \mathbf{y}$$
$$\uparrow$$
$$W$$

# Examples of computational units in NN:
## Output and its derivatives w.r.t. inputs and parameters

- Linear unit:

  - Output: $\mathbf{z} = \boldsymbol{W}^{\top}\mathbf{x}$

  - Derivatives: $\partial\mathbf{z}/\partial\mathbf{x} = \boldsymbol{W}, \quad \partial\mathbf{z}/\partial\boldsymbol{W} = \mathbf{x}$

- Sigmoid unit (Logistic unit):

  - Output: $\mathbf{z} = \sigma(\mathbf{x})$ ($\sigma$ is element-wise application of sigmoid function)

  - Derivatives: $\dfrac{\partial\mathbf{z}}{\partial\mathbf{x}} = \sigma(\mathbf{x})(1 - \sigma(\mathbf{x}))$ (No parameter derivative)

- ReLU unit (Rectangular Linear Unit):

  - Output: $\mathbf{z} = \max\{\mathbf{0}, \mathbf{x}\}$ (element-wise max)

  - Derivatives: $\dfrac{\partial\mathbf{z}}{\partial\mathbf{x}} = \min\{\mathbf{0}, \mathbf{1}\}$ (No parameter derivative)

# Automatic differentiation over computational graph: Forward path

- All units are differentiable w.r.t. inputs & parameters

- Traversals on backward paths give arbitrary derivatives



$$\frac{\partial \mathbf{z}_3}{\partial \mathbf{z}_1} = \sigma(\mathbf{z}_1)\big(1 - \sigma(\mathbf{z}_1)\big)$$

$$\frac{\partial \mathbf{z}_1}{\partial \mathbf{x}^{(i)}} = W_1$$

$$\frac{\partial \mathbf{z}_1}{\partial W_1} = \mathbf{x}$$

$$\frac{\partial \mathbf{z}_5}{\partial \mathbf{z}_3} = \operatorname{diag}(\mathbf{z}_4)$$

$$\frac{\partial y}{\partial z_6} = \sigma(z_6)\big(1 - \sigma(z_6)\big)$$

$$\frac{\partial \mathbf{z}_6}{\partial \mathbf{z}_5} = \mathbf{w}$$

$$\frac{\partial \mathbf{z}_4}{\partial \mathbf{z}_2} = \min\{\mathbf{0}, \mathbf{1}\}$$

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{x}^{(i)}} = W_2$$

$$\frac{\partial \mathbf{z}_5}{\partial \mathbf{z}_4} = \operatorname{diag}(\mathbf{z}_3)$$

$$\frac{\partial \mathbf{z}_6}{\partial \mathbf{w}} = \mathbf{z}_5$$

$$\frac{\partial L}{\partial y} = \frac{\delta(y^{(i)} = 1)}{y} + \frac{\delta(y^{(i)} = -1)}{1 - y}$$

$$\frac{\partial \mathbf{z}_2}{\partial W_2} = \mathbf{x}$$

$$L = CE(y_7, y^{(i)})$$

# Automatic differentiation over computational graph: Backward path

- Traversals on backward paths give arbitrary derivatives



$$\frac{\partial \mathbf{z}_3}{\partial \mathbf{z}_1} = \sigma(\mathbf{z}_1)\big(1 - \sigma(\mathbf{z}_1)\big)$$

$$\frac{\partial \mathbf{z}_1}{\partial \mathbf{x}^{(i)}} = W_1$$

$f_1$

$$\frac{\partial \mathbf{z}_1}{\partial W_1} = \mathbf{x}$$

$W_1$

$\sigma$

$$\frac{\partial \mathbf{z}_5}{\partial \mathbf{z}_3} = \mathrm{diag}(\mathbf{z}_4)$$

$$\frac{\partial y}{\partial z_6} = \sigma(z_6)\big(1 - \sigma(z_6)\big)$$

$\mathbf{x}^{(i)}$

$\circ$

$f_3$

$\sigma$

$$\frac{\partial z_6}{\partial \mathbf{z}_5} = \mathbf{w}$$

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{x}^{(i)}} = W_2$$

$$\frac{\partial \mathbf{z}_4}{\partial \mathbf{z}_2} = \min\{\mathbf{0}, \mathbf{1}\}$$

$f_2$

ReLU

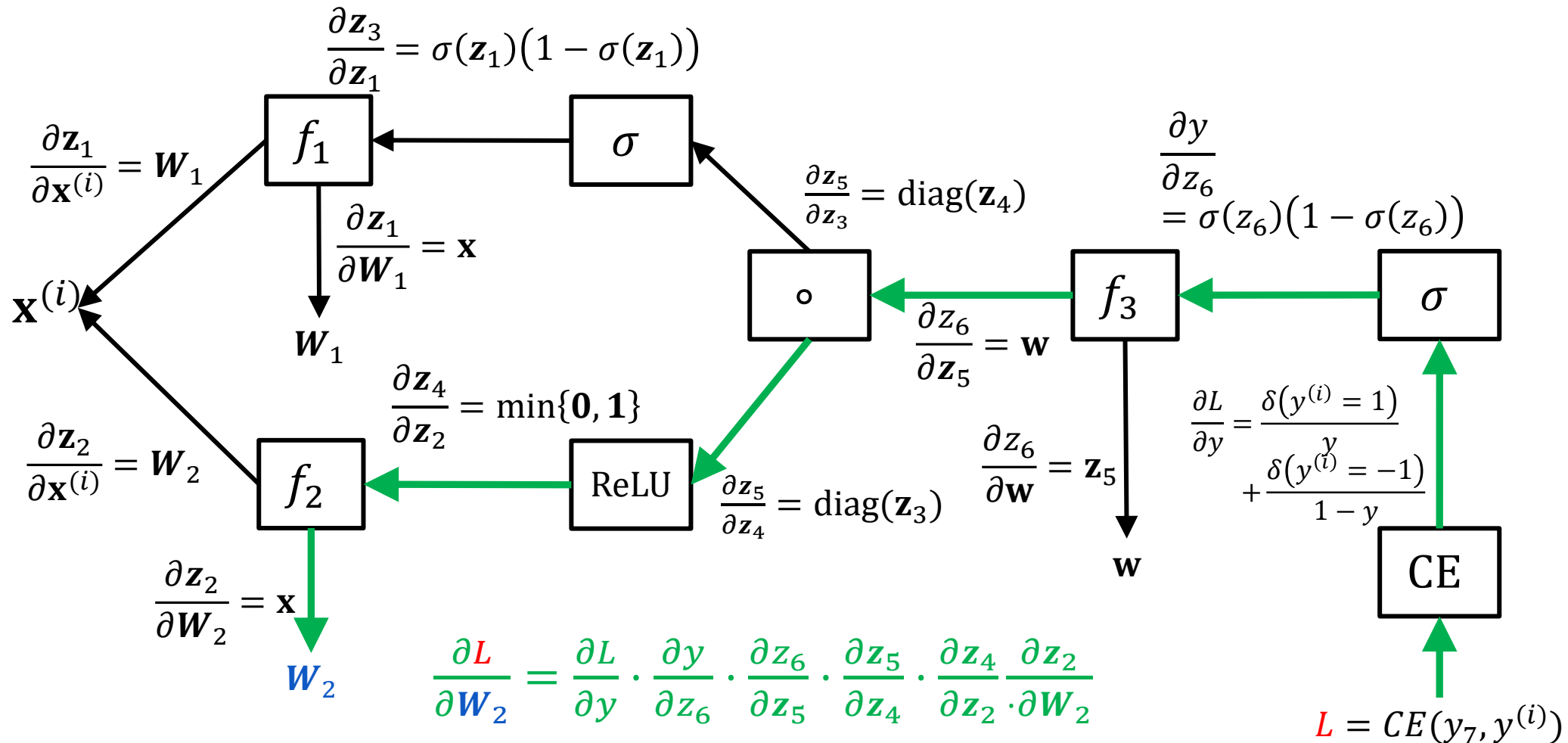$$\frac{\partial \mathbf{z}_5}{\partial \mathbf{z}_4} = \mathrm{diag}(\mathbf{z}_3)$$

$$\frac{\partial z_6}{\partial \mathbf{w}} = \mathbf{z}_5$$

$$\frac{\partial L}{\partial y} = \frac{\delta(y^{(i)} = 1)}{y} + \frac{\delta(y^{(i)} = -1)}{1 - y}$$

$$\frac{\partial \mathbf{z}_2}{\partial W_2} = \mathbf{x}$$

$W_2$

$\mathbf{w}$

CE

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z_6} \cdot \frac{\partial z_6}{\partial \mathbf{z}_5} \cdot \frac{\partial \mathbf{z}_5}{\partial \mathbf{z}_4} \cdot \frac{\partial \mathbf{z}_4}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial W_2}$$

$$L = CE(y_7, y^{(i)})$$

# Various deep neural network structure:
## Flexible modeling by combining units

- NN allows intuitive and flexible modeling by combining various types of units like "*LEGO®*" blocks to form complex networks



BN: Batch Normalization
ReLU: Rectified Linear Unit

- Error back propagation can be left to DL frameworks such as PyTorch

- Various task dependent neural networks have been proposed:

  - Images: Convolutional Neural Network (CNN); Vision Transformer (ViT)

  - Languages: Recurrent Neural Network (RNN), Transformer,

  - Graphs: Graph neural networks (GNN)
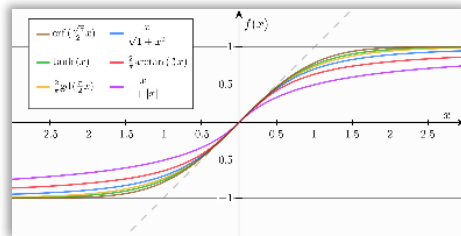
# Training Deep Networks

# Gradient vanishing problem:
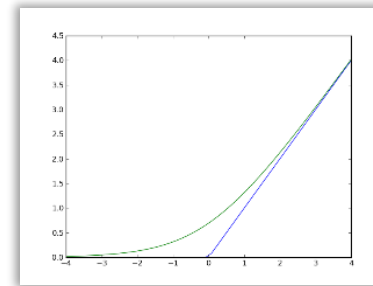## A major obstacle in training deep networks

- You can design and train arbitrary large networks, but in practice, we face the problem that training does not go well..

- Gradient vanishing problem: the derivative becomes weaker and weaker in the process of error back propagation

  - In the previous example, $\dfrac{\partial L}{\partial \boldsymbol{W}_2} = \dfrac{\partial L}{\partial y} \cdot \dfrac{\partial y}{\partial z_6} \cdot \dfrac{\partial z_6}{\partial z_5} \cdot \dfrac{\partial z_5}{\partial \boldsymbol{z}_4} \cdot \dfrac{\partial \boldsymbol{z}_4}{\partial \boldsymbol{z}_2} \dfrac{\partial \boldsymbol{z}_2}{\cdot \partial \boldsymbol{W}_2}$

  - Repeated multiplication of small derivative values results in progressively smaller overall derivative values

- In deep networks, the gradient tends to be small, and hence gradient descent optimization does not proceed

# How to deal with gradient vanishing problem: ReLU and skip connection

- Gradient of logistic function becomes small away from the origin

- Activation functions other than logistic function

  - Rectangular linear unit: $\max(0, x)$
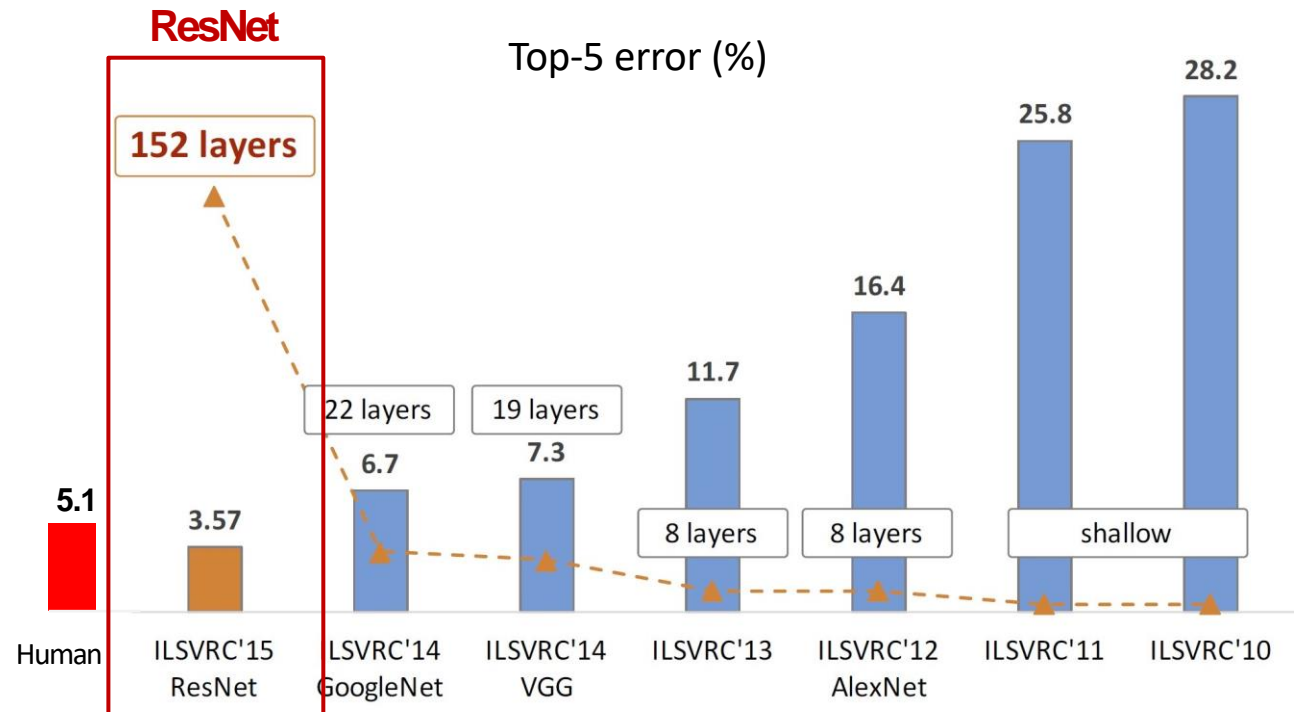
    - Its derivative is 1 (or 0)
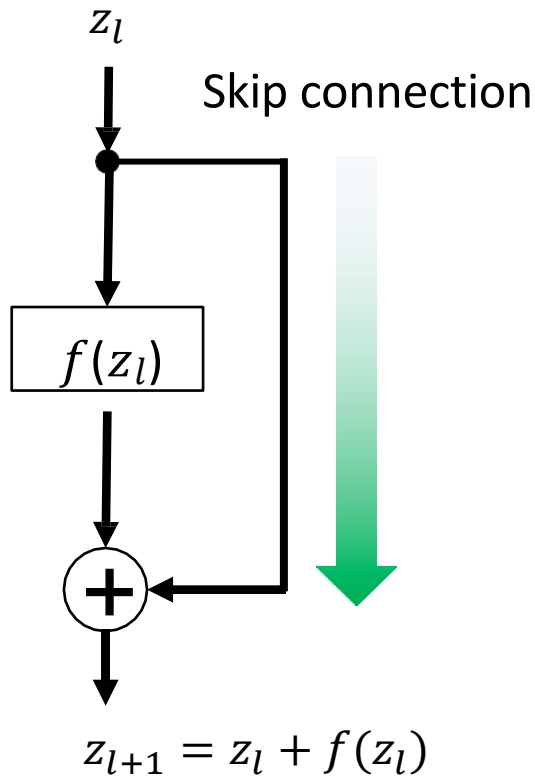


Logistic function

ReLU

    - Softplus: $\ln(1 + e^x)$, LeakyReLU, …. and others

# How to deal with gradient vanishing problem:
## ReLU and skip connection

- Skip connection: $z_{l+1} = z_l + f(z_l)$

- Gradient is propagated directly without decay



$z_l$

Skip connection

$f(z_l)$

$+$

$z_{l+1} = z_l + f(z_l)$

**ResNet**

Top-5 error (%)

152 layers

22 layers

19 layers

8 layers

8 layers

shallow

5.1

3.57

6.7

7.3

11.7

16.4

25.8

28.2

Human

ILSVRC'15 ResNet

ILSVRC'14 GoogleNet

ILSVRC'14 VGG

ILSVRC'13

ILSVRC'12 AlexNet

ILSVRC'11

ILSVRC'10

# Summary:
# Neural networks

1. Neural network: (Very roughly speaking) stacked logistic regression

2. Training neural networks: gradient method (SGD, mini-batch, …)

3. Computational graphs and automatic differentiation (error back propagation): gradients can be computed efficiently and systematically on a computational graph consisting of simple differentiable units

4. Training deep learning models: Dealing with the Gradient vanishing problem with ReLU, skip connection, …