# Inferencing Issues:
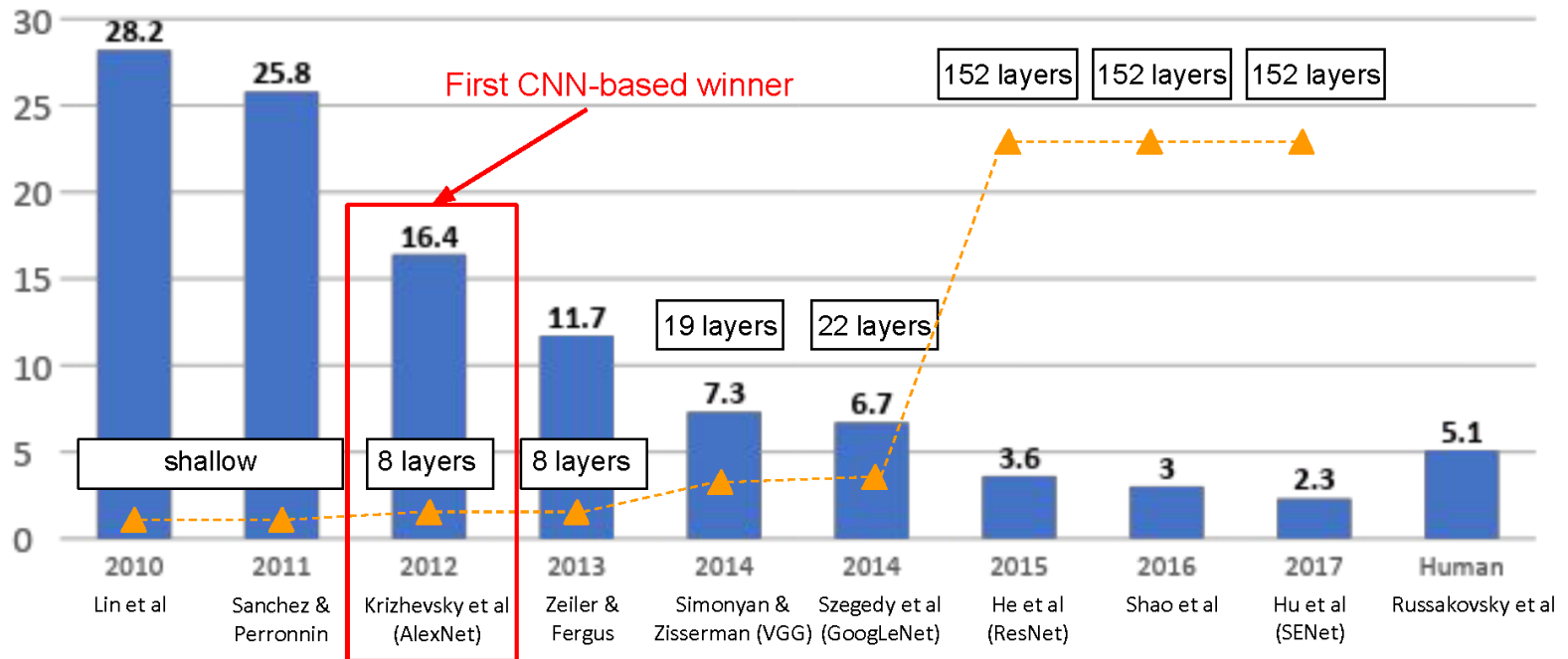# Overfitting and Learning Dilemma

國立政治大學 資訊管理學系

蔡瑞煌 特聘教授

# AI applications

- Training phase: (training) data + AI model + algorithm & code + setting of network & hyperparameters → AI model/AI system

- Inferencing phase:  performance is obtained from model((test) data)

- Goals of training are reasonable inferencing

# CNN models

## ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

idea/concept of learning →

a learning algorithm →
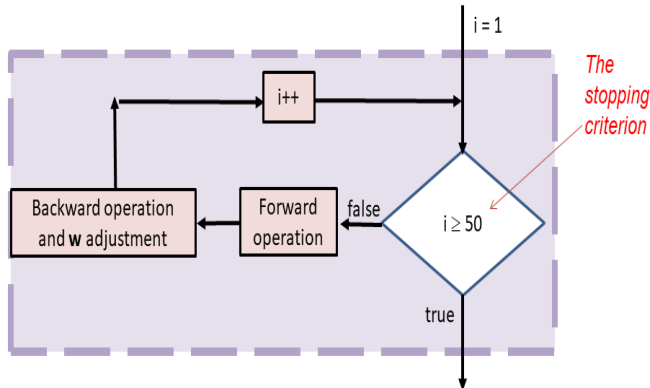
codes →

an AI model/system

# Algorithm

(Algorithm - Wikipedia)

- In mathematics and computer science, an **algorithm** is a finite sequence of well-defined, computer-implementable instructions, typically to solve a class of problems or to perform a computation.

✓ Initial state/system
✓ Finite steps/blocks/modules
✓ Sequential order ($\rightarrow$)
✓ Loop (for 迴圈； iteration/epoch)
✓ Goal
✓ Stopping criteria

Where we are now...

# TensorFlow: Loss

```
N, D, H = 64, 1000, 100

x = tf.convert_to_tensor(np.random.randn(N, D), np.float32)
y = tf.convert_to_tensor(np.random.randn(N, D), np.float32)
w1 = tf.Variable(tf.random.uniform((D, H)))   # weights
w2 = tf.Variable(tf.random.uniform((H, D)))   # weights

optimizer = tf.optimizers.SGD(1e-6)

for t in range(50):
    with tf.GradientTape() as tape:
        h = tf.maximum(tf.matmul(x, w1), 0)
        y_pred = tf.matmul(h, w2)
        diff = y_pred - y
        loss = tf.losses.MeanSquaredError()(y_pred, y)
    gradients = tape.gradient(loss, [w1, w2])
    optimizer.apply_gradients(zip(gradients, [w1, w2]))
```
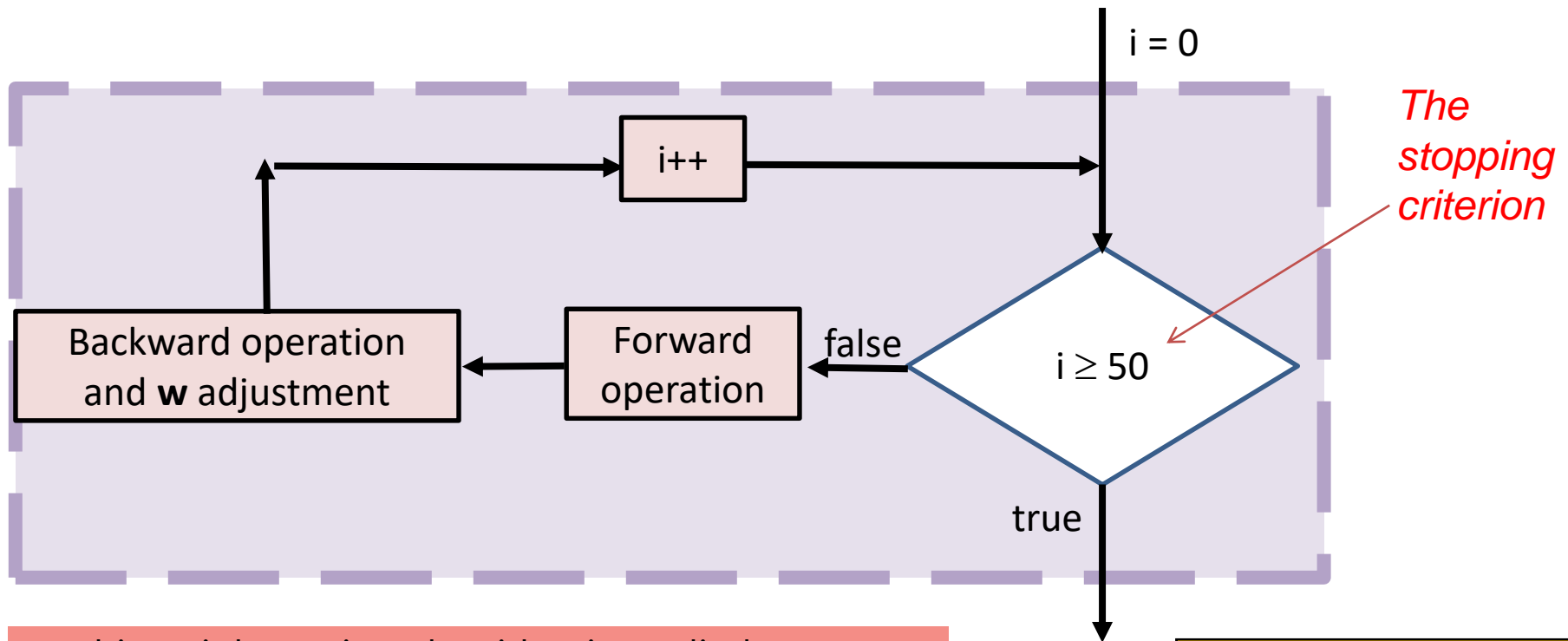
Use predefined loss functions

The flowchart form of algorithm



i = 1

The stopping criterion

Backward operation and **w** adjustment

Forward operation

false

i ≥ 50

true

i++

Weight-tuning module_EU

# The flowchart of weight-tuning algorithms for 2-layer neural networks in CS231n

i = 0

*The stopping criterion*



i++

Backward operation and **w** adjustment
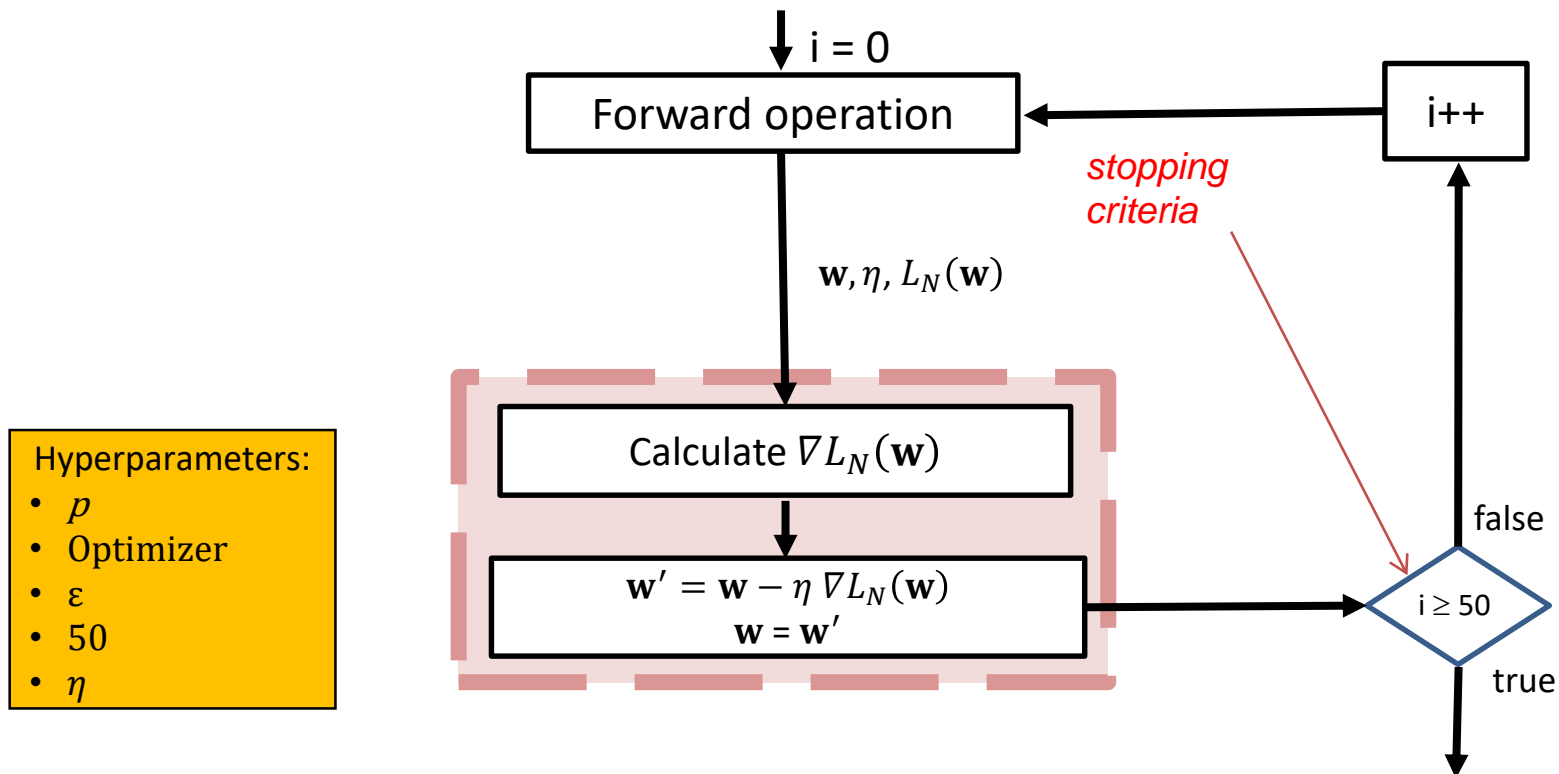
Forward operation

false

i ≥ 50

true

- This weight-tuning algorithm is applied to many kinds of Neural Networks, including 2-layer neural networks, CNN, RNN, reinforcement learning, GAN, BERT, and so on.
- The weight-tuning process stops when the stopping criterion is satisfied.

Hyperparameters:
- $p$
- Optimizer: SGD
- Epoch upper bound: 50
- Learning rate: 1e-6

7

# The flowchart of weight-tuning module_EU



Hyperparameters:
- $p$
- Optimizer
- $\varepsilon$
- 50
- $\eta$

i = 0

Forward operation

i++

stopping criteria

$\mathbf{w}, \eta, L_N(\mathbf{w})$

Calculate $\nabla L_N(\mathbf{w})$

$\mathbf{w}' = \mathbf{w} - \eta \, \nabla L_N(\mathbf{w})$
$\mathbf{w} = \mathbf{w}'$

$i \geq 50$

false

true

8

# The learning goals
# (also the stopping criteria for the learning)

The learning process should stop when

1. ~~$L_N(\mathbf{w}) = 0$~~

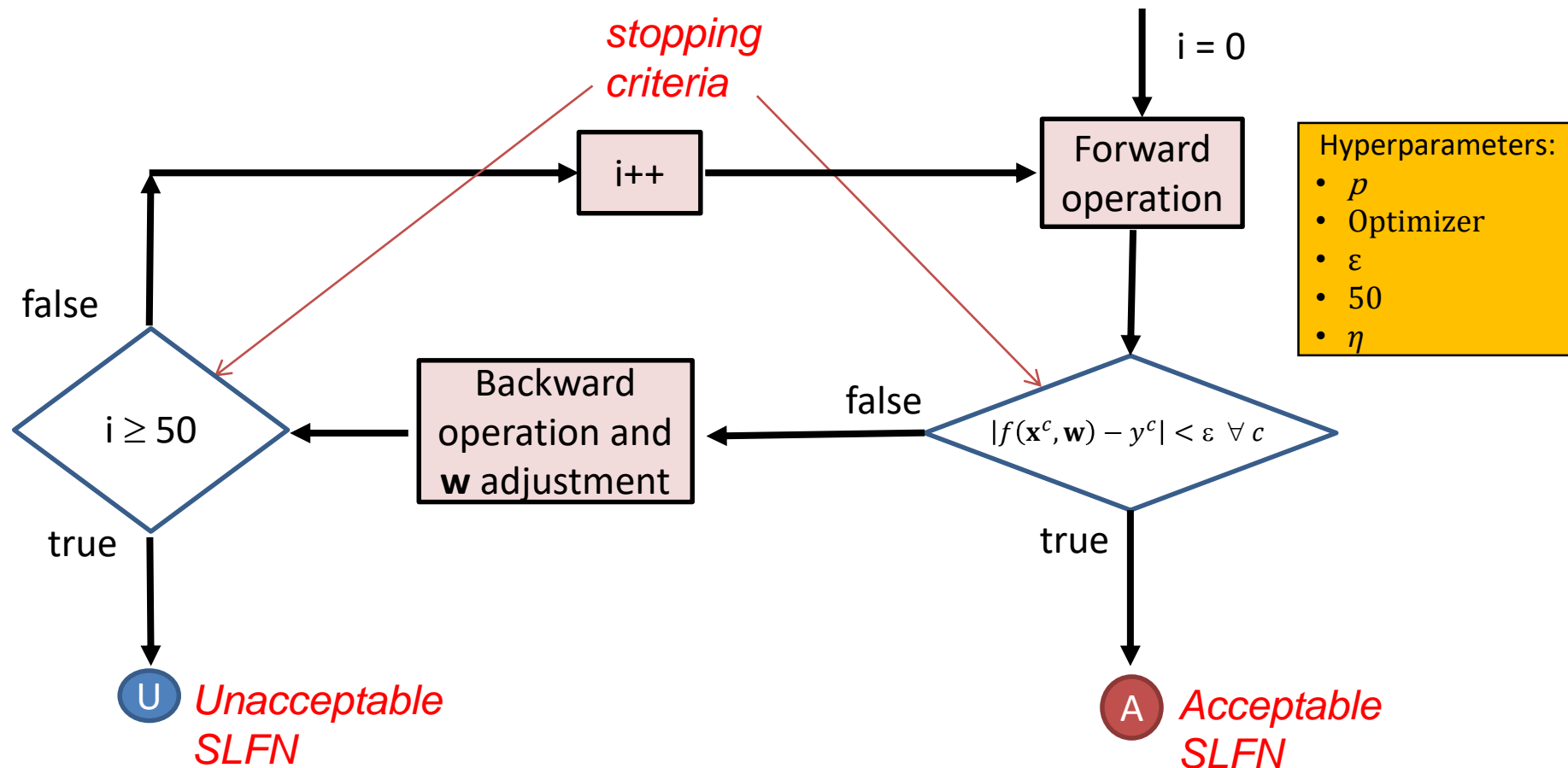2. a tiny $L_N(\mathbf{w})$ value

3. $|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \ c \ \text{with } \varepsilon \text{ being tiny}$

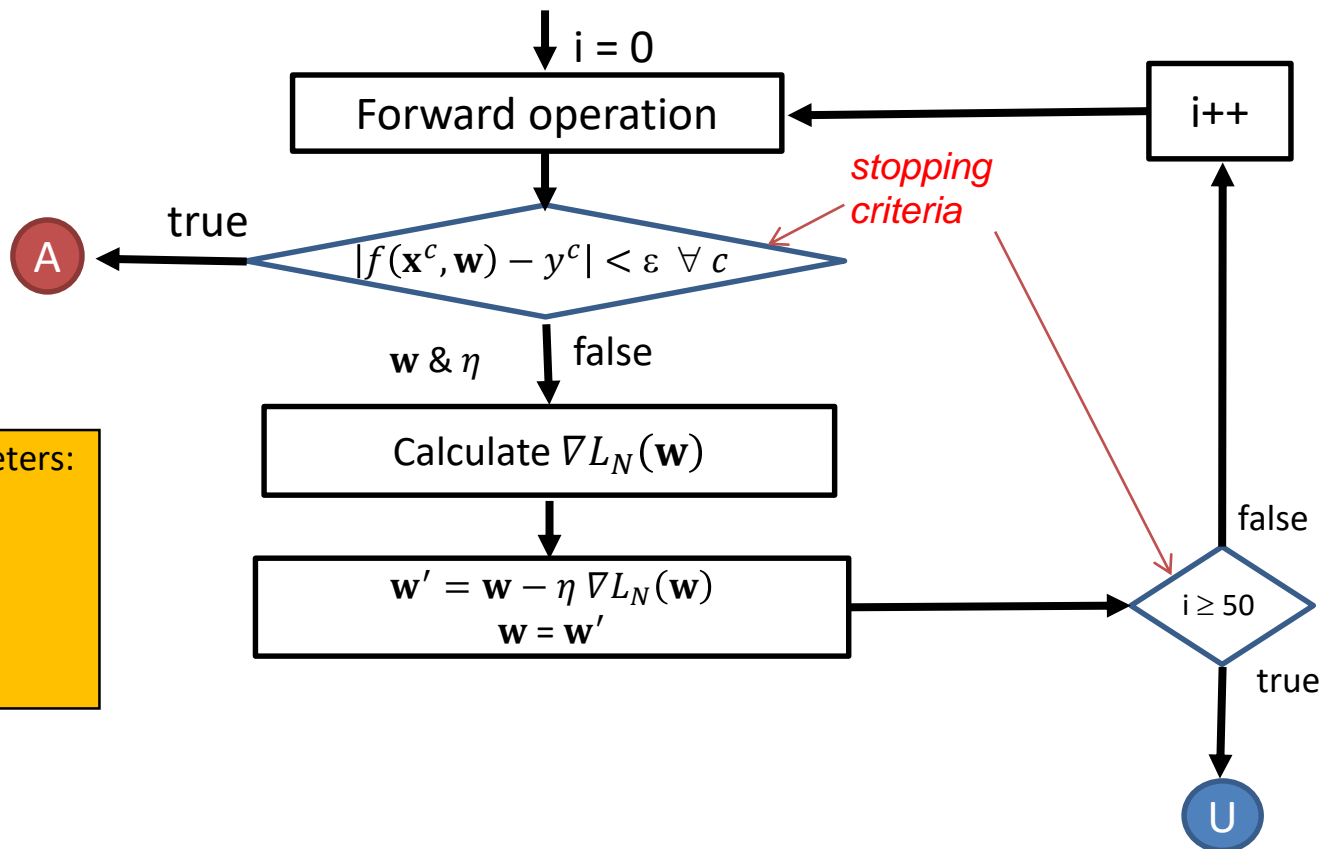$$L_N(\mathbf{w}) \equiv \frac{1}{N} \sum_{c=1}^{N} (f(\mathbf{x}^c, \mathbf{w}) - y^c)^2$$

- Each reasonable learning goal can be used as a stopping criterion.
- Different stopping criterion results in different length of training time and different model.

Where we are now...

# The flowchart of weight-tuning algorithm including two stopping criteria that indicate either an unacceptable SLFN or an acceptable SLFN

*stopping criteria*

i = 0

i++

Forward operation

Hyperparameters:
- $p$
- Optimizer
- $\varepsilon$
- 50
- $\eta$

false

i ≥ 50

Backward operation and **w** adjustment

false

$|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \, c$

true

true

U *Unacceptable SLFN*

A *Acceptable SLFN*

10

# The flowchart of <span style="color:red">weight-tuning</span> module_EU_LG



Hyperparameters:
- $p$
- Optimizer
- $\varepsilon$
- 50
- $\eta$

$i = 0$

Forward operation

$i{+}{+}$

$|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \;\; \forall \, c$

*stopping criteria*

true → A

$\mathbf{w} \,\&\, \eta$  false

Calculate $\nabla L_N(\mathbf{w})$

$\mathbf{w}' = \mathbf{w} - \eta \, \nabla L_N(\mathbf{w})$
$\mathbf{w} = \mathbf{w}'$

$i \geq 50$

false

true

U

# Learning dilemma of gradient-descent-based learning



loss

Very slow at the plateau and neighborhoods of attractors

Stuck at saddle points

Stuck at local minima

$$\nabla L(w) \approx 0$$

$$\nabla L(w) = 0$$

$$\nabla L(w) = 0$$

**w** space

# Extra stopping criteria for the learning (not the learning goals)

$\|\nabla_{\mathbf{w}} L_N(\mathbf{w})\|$ is the length of $\nabla_{\mathbf{w}} L_N(\mathbf{w})$.
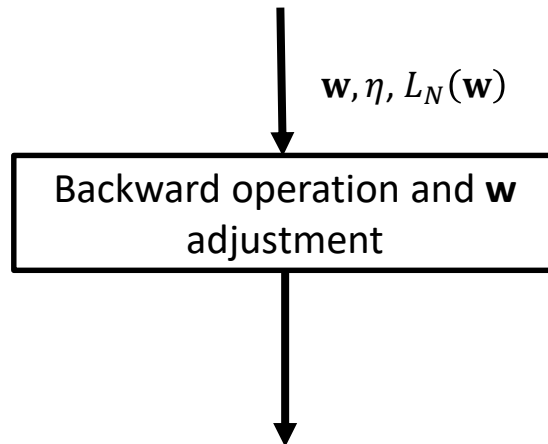
1.  ~~The learning process should stop when $\|\nabla_{\mathbf{w}} L_N(\mathbf{w})\| = 0$ but a tiny $L_N(\mathbf{w})$ value cannot be accomplished.~~
2.  The learning process should stop when $\|\nabla_{\mathbf{w}} L_N(\mathbf{w})\|$ is tiny but a tiny $L_N(\mathbf{w})$ value cannot be accomplished.
3.  The learning process should stop when $\eta$ (the learning rate) is tiny but a tiny $L_N(\mathbf{w})$ value cannot be accomplished.

The neighborhood of undesired attractors, where $\|\nabla_{\mathbf{w}} L_N(\mathbf{w})\| \approx 0$ but a tiny $L_N(\mathbf{w})$ value cannot be accomplished:
  a) the local minimum, the saddle point, or the plateau.
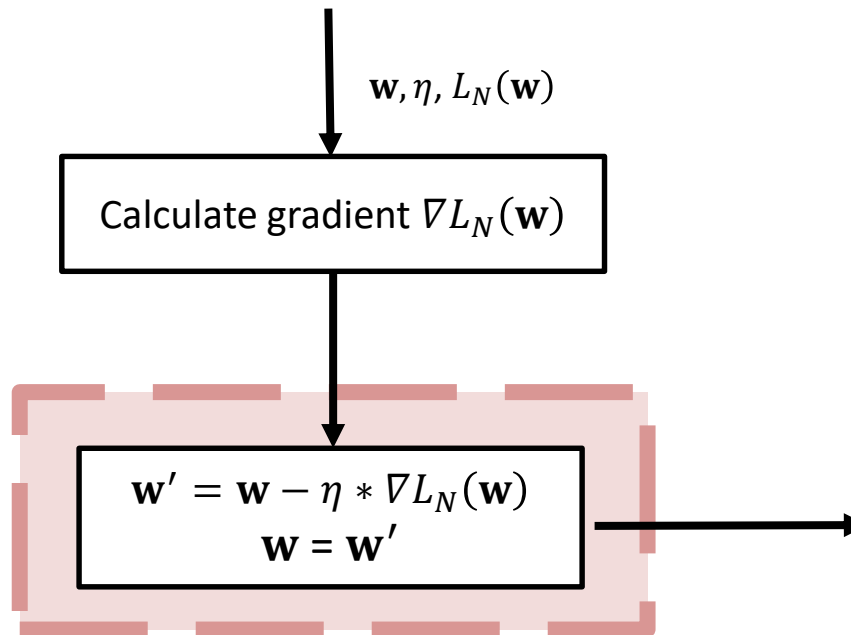  b) the global minimum of the defective network architecture.
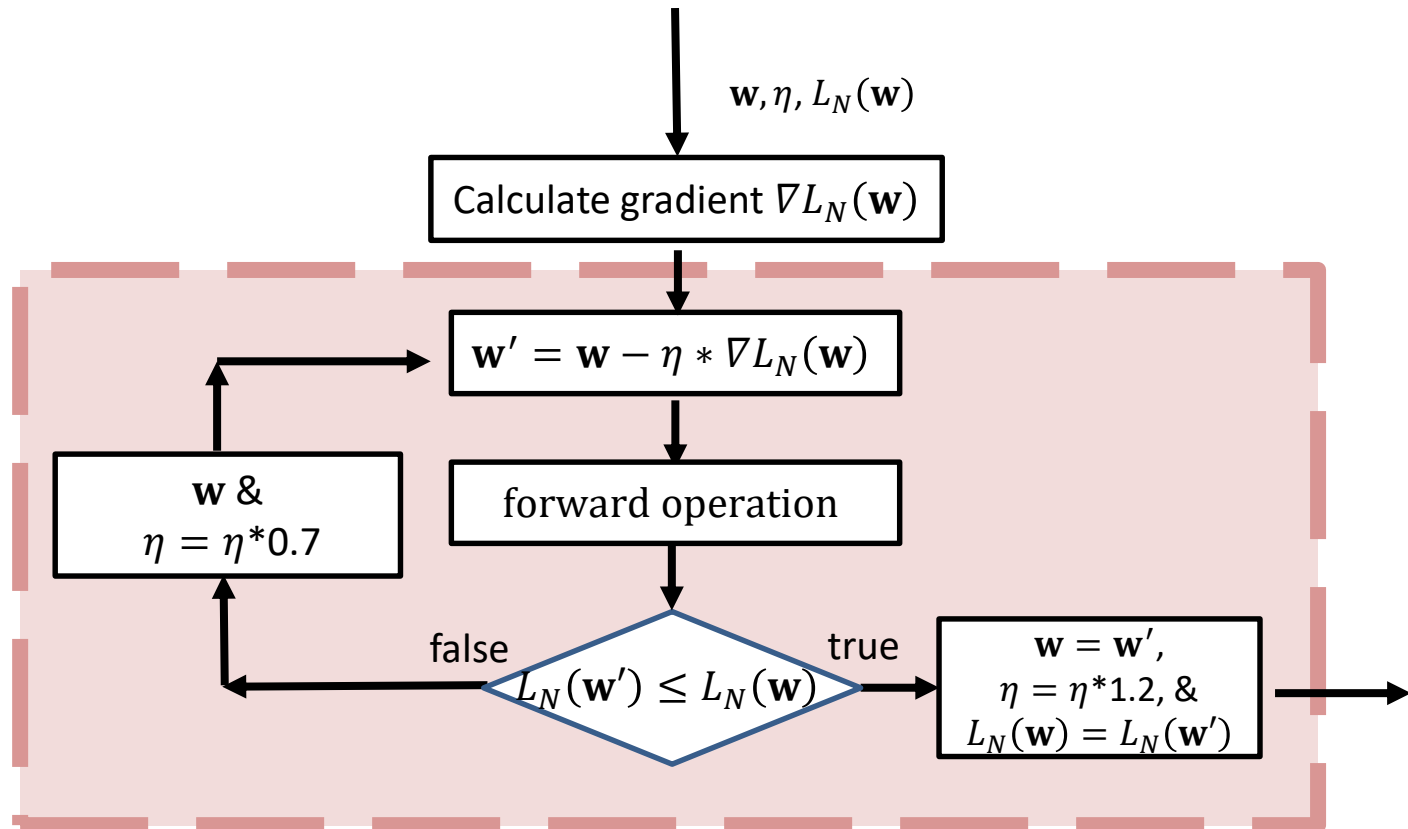
# Module of backward operation and **w** adjustment

$$\mathbf{w}, \eta, L_N(\mathbf{w})$$

Backward operation and **w** adjustment

# Module of backward operation and **w** adjustment

$$\mathbf{w}, \eta, L_N(\mathbf{w})$$

Calculate gradient $\nabla L_N(\mathbf{w})$

$$\mathbf{w}' = \mathbf{w} - \eta * \nabla L_N(\mathbf{w})$$
$$\mathbf{w} = \mathbf{w}'$$

# Module of backward operation and **w** adjustment with the adaptable learning rate arrangement

$$\mathbf{w}, \eta, L_N(\mathbf{w})$$

Calculate gradient $\nabla L_N(\mathbf{w})$

$$\mathbf{w}' = \mathbf{w} - \eta * \nabla L_N(\mathbf{w})$$

**w** &
$\eta = \eta * 0.7$

forward operation

false $\qquad$ $L_N(\mathbf{w}') \leq L_N(\mathbf{w})$ $\qquad$ true

$\mathbf{w} = \mathbf{w}'$,
$\eta = \eta * 1.2$, &
$L_N(\mathbf{w}) = L_N(\mathbf{w}')$

16

# The flowchart of weight-tuning module_EU_LG

**Hyperparameters:**
- $p$
- Optimizer
- $\varepsilon$
- 50
- 1.2 & 0.7

i = 0

Forward operation

i++

**stopping criteria**

true

A

$\varepsilon \geq L_N(\mathbf{w})$

$\mathbf{w}, \eta, L_N(\mathbf{w})$  false

Calculate $\nabla L_N(\mathbf{w})$

$\mathbf{w}' = \mathbf{w} - \eta * \nabla L_N(\mathbf{w})$

false

$i \geq 50$

forward operation

$\mathbf{w}$ &
$\eta = \eta * 0.7$

true

U

false    $L_N(\mathbf{w}') \leq L_N(\mathbf{w})$    true

$\mathbf{w} = \mathbf{w}'$,
$\eta = \eta * 1.2$, &
$L_N(\mathbf{w}) = L_N(\mathbf{w}')$

17

# The effect of adaptable learning rate



Loss / Mean Squared Error

# The flowchart of weight-tuning module_EU_LG_UA



$i = 0$

Forward operation

*stopping criteria*

A

true

$|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \ c$

**Hyperparameters:**
- $p$
- Optimizer
- $\varepsilon \ \& \ \varepsilon_1$
- 1.2 & 0.7
- 50

$\mathbf{w}, \eta, L_N(\mathbf{w})$  false

Calculate $\nabla L_N(\mathbf{w})$

i++

$\mathbf{w}' = \mathbf{w} - \eta \ \nabla L_N(\mathbf{w})$

The mechanism of identifying the neighborhood of an undesired attractor

Forward operation

$\mathbf{w} \ \& \ 0.7\eta \to \eta$

$L_N(\mathbf{w}') < L_N(\mathbf{w})$  true

$\mathbf{w} = \mathbf{w}'$,
$\eta = \eta*1.2$, &
$L_N(\mathbf{w}) = L_N(\mathbf{w}')$

$i \geq 50$  false

false

true

U

$\eta > \varepsilon_1$  false

true

19

# The flowchart of BP learning algorithm



Forward operation

*Capable of stopping?*

Backward operation and **w** adjustment

false $|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \ c$

true

*???* (A) *Acceptable SLFN*

Hyperparameters:
- $p$
- Optimizer
- $\varepsilon$
- $\eta$

Q: Can this the learning process stop through satisfying the stopping criterion?
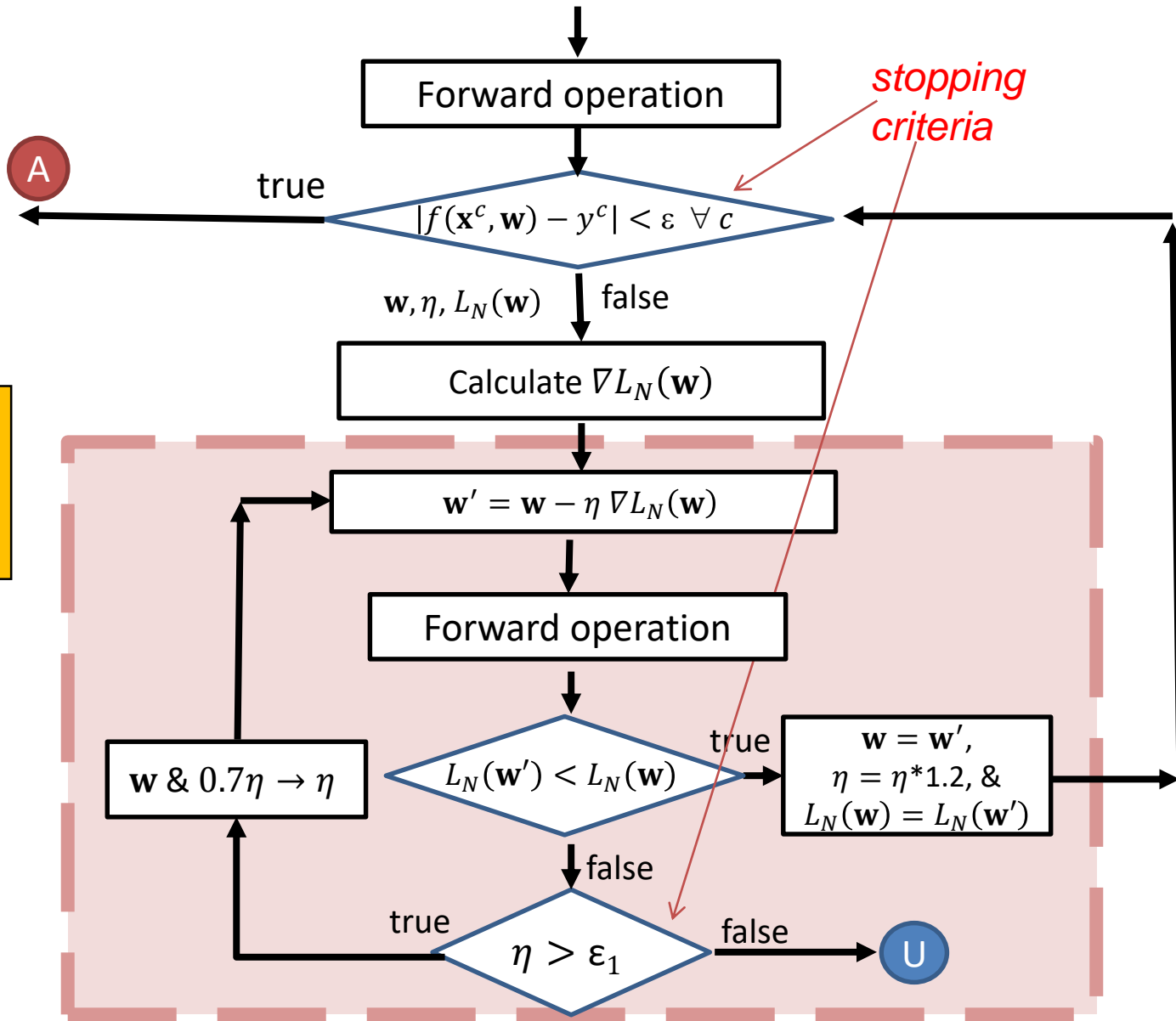A: May not be! So, this stopping criterion is not good and thus this algorithm is not good.

# The flowchart of BP learning algorithm



Forward operation

true $\qquad |f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \;\; \forall\, c$

A

$\mathbf{w}, \eta, L_N(\mathbf{w})$ \qquad false

Calculate $\nabla L_N(\mathbf{w})$

$$\mathbf{w}' = \mathbf{w} - \eta\, \nabla L_N(\mathbf{w})$$
$$\mathbf{w} = \mathbf{w}'$$

Hyperparameters:
- $p$
- Optimizer
- $\varepsilon$
- $\eta$

# The flowchart of weight-tuning module_LG_UA



Forward operation

*stopping criteria*

$|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall c$

A

true

$\mathbf{w}, \eta, L_N(\mathbf{w})$   false

Calculate $\nabla L_N(\mathbf{w})$

Hyperparameters:
- $p$
- Optimizer
- $\varepsilon \ \& \ \varepsilon_1$
- 1.2 & 0.7

$\mathbf{w}' = \mathbf{w} - \eta \, \nabla L_N(\mathbf{w})$

Forward operation

$\mathbf{w} \ \& \ 0.7\eta \to \eta$

$L_N(\mathbf{w}') < L_N(\mathbf{w})$

true

$\mathbf{w} = \mathbf{w}'$,
$\eta = \eta * 1.2$, &
$L_N(\mathbf{w}) = L_N(\mathbf{w}')$

false

true

$\eta > \varepsilon_1$

false

U

22

# Performance differences amongst weight-tuning modules

- There are four weight-tuning modules
  - ✓ the weight-tuning module_EU
  - ✓ the weight-tuning module_EU_LG
  - ✓ the weight-tuning module_EU_LG_UA
  - ✓ the weight-tuning module_LG_UA
- What are the performance differences amongst these weight-tuning modules?

# Performance differences amongst weight-tuning modules

- There are four weight-tuning modules
  - ✓ the weight-tuning module_EU

    The simplest and the learning time length is expected

  - ✓ the weight-tuning module_EU_LG

    Shorter learning time length than the weight-tuning module_EU

  - ✓ the weight-tuning module_EU_LG_UA

    The learning time length may be longer than the weight-tuning module_EU_LG

  - ✓ the weight-tuning module_LG_UA

    The learning time length is not the issue

# Homework #2

- Rewrite the code you have for HW #1 (the weight-tuning module_EU referring to page 8) into the code of the weight-tuning module_EU_LG referring to page 11.
- Rewrite the code you have for HW #1 (the weight-tuning module_EU) into the code of the weight-tuning module_EU_LG_UA referring to page 19.
- Rewrite the code you have for HW #1 (the weight-tuning module_EU) into the code of the weight-tuning module_LG_UA referring to page 22.
- Once you have the code (regardless of which framework you choose above), you will apply the code to learn the train_all_0.csv dataset given in the LINE group.
- The training and test dataset is 80%/20%.
- The performance comparison benchmark is the code of the weight-tuning module_EU.
- Note that the output nodes of your SLFN may be one or two.

# The supervised learning problems: Regression and Classification

**Regression**

What is the temperature going to be tomorrow?

PREDICTION

84°

Fahrenheit °F
-50 -40 -30 -20 -10 0 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150 160 170 180 190 200 210 220 230

**Classification**

Will it be Cold or Hot tomorrow?

PREDICTION

COLD

HOT

Fahrenheit °F
-50 -40 -30 -20 -10 0 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150 160 170 180 190 200 210 220 230

# Stopping criteria (also the learning goals) for regression problems

The learning process should stop when

1. ~~$L_N(\mathbf{w}) = 0$~~

2. a tiny $L_N(\mathbf{w})$ value

3. $|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \ c \ with \ \varepsilon \ being \ tiny$

one output node

$$L_N(\mathbf{w}) \equiv \frac{1}{N}\sum_{c=1}^{N}(f(\mathbf{x}^c, \mathbf{w}) - y^c)^2$$

- Each reasonable learning goal can be used as a stopping criterion.
- Different stopping criterion results in different length of training time and different model.

# The regression applications

## The learning mechanism

$$|f(\mathbf{x}^c, \mathbf{w}) - y^c| \leq \varepsilon \ \forall \ c \in \mathbf{I}$$



$f(\mathbf{x}^c, \mathbf{w})$

$y^c - \varepsilon$ $y^c + \varepsilon$

$y^c$

## The inferencing mechanism



not $y^c$

$f(\mathbf{x}^c, \mathbf{w})$

$y^c - \varepsilon$ $y^c + \varepsilon$

$y^c$

# The three-class classification problems

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Algebraic Viewpoint

Visual Viewpoint

$$f(x,W) = Wx$$

29

# The three-class classification problems

Softmax vs. SVM

matrix multiply + bias offset

hinge loss (SVM)

max(0, -2.85 - 0.28 + 1) +
max(0, 0.86 - 0.28 + 1)
=
**1.58**

cross-entropy loss (Softmax)

- log(0.353)
=
**0.452**

$W$     $x_i$     $b$

$y_i$ | 2

30

# Classification Applications Design (*y* label)

**Output value: real number**

*SLFN with* one output node and linear arrangement

疲倦
- ✓ 無: 46位
- ✓ 輕度:346 位
- ✓ 中至重度:294位

Learning phase:
*y* (i.e., target output):
- ✓ 無: 0
- ✓ 輕度: 5
- ✓ 中至重度: 10

Inferencing phase:
*f* (i.e., actual output):
- ✓ [-2.5, 2.5) → 無
- ✓ [2.5, 7.5) → 輕度
- ✓ [7.5, 12.5) →中至重度
- ✓ (-∞, -2.5) OR [12.5, ∞) → unknown

**Output value: binary number**

*SLFN with three* output nodes and softmax arrangement

疲倦
- ✓ 無: 46位
- ✓ 輕度:346 位
- ✓ 中至重度:294位

Learning phase:
*y* (i.e., target output):
- ✓ 無: (1, 0, 0)
- ✓ 輕度: (0, 1, 0)
- ✓ 中至重度: (0, 0, 1)

Inferencing phase:
*f* (i.e., actual output):
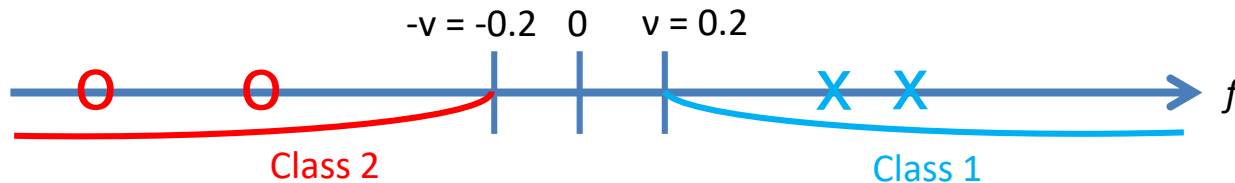- ✓ (1, 0, 0) → 無
- ✓ (0, 1, 0) → 輕度
- ✓ (0, 0, 1) →中至重度

# Stopping criteria (also the learning goals) for two-class classification problems dealt with the SLFN with one output node whose output values are real numbers

- Two-class classification problems with $\mathbf{I} \equiv \mathbf{I}_1 \cup \mathbf{I}_2$, where $\mathbf{I}_1$ and $\mathbf{I}_2$ are the sets of indices of given cases in classes 1 and 2. Furthermore, $y^c$ is the target of the $c^{\text{th}}$ case, with 1 and 0 being the targets of classes 1 and 2

  one output node

- When the SLFN with only one output node whose output value is real number, the stopping criteria may be as follows:

1. $|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \, c$

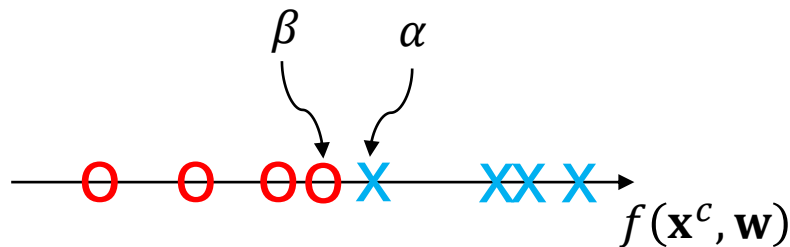2. $f(\mathbf{x}^c, \mathbf{w}) > v \ \forall \, c \in \mathbf{I}_1$ and $f(\mathbf{x}^c, \mathbf{w}) \leq -v \ \forall \, c \in \mathbf{I}_2,$ with $1 > v > 0$

3. $\alpha \equiv \min_{c \epsilon \mathbf{I}_1} f(\mathbf{x}^c, \mathbf{w}) > \beta \equiv \max_{c \epsilon \mathbf{I}_2} f(\mathbf{x}^c, \mathbf{w})$
   (Linearly seperating condition, $LSC$)

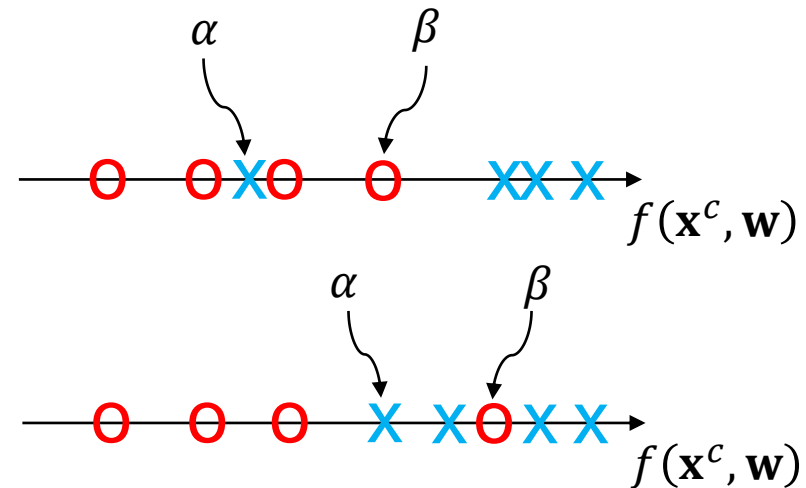   Different stopping criterion results in different length of training time and different model.

# Stopping criteria (also the learning goals) for two-class classification problems dealt with the SLFN with one output node whose output values are real numbers

$y^c = 1 \ \forall \ c \in \mathbf{I}_1; \ y^c = 0 \ \forall \ c \in \mathbf{I}_2$    X :  $f(\mathbf{x}^c, \mathbf{w}), \ \forall \ c \in \mathbf{I}_1$    O :  $f(\mathbf{x}^c, \mathbf{w}), \ \forall \ c \in \mathbf{I}_2$
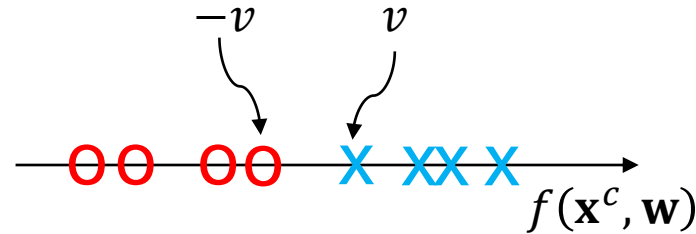
● $|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \ c$



● $f(\mathbf{x}^c, \mathbf{w}) \geq v \ \forall \ c \in \mathbf{I}_1$ and $f(\mathbf{x}^c, \mathbf{w}) \leq -v \ \forall \ c \in \mathbf{I}_2$ with $1 > v > 0$.

# Stopping criteria (also the learning goals) for two-class classification problems dealt with the SLFN with one output node whose output values are real numbers

● The LSC regarding $\{f(\mathbf{x}^c, \mathbf{w}) \ \forall \ c \in \mathbf{I}\}$    (Tsaih, 1993)

$y^c = 1 \ \forall \ c \in \mathbf{I}_1; \ y^c = 0 \ \forall \ c \in \mathbf{I}_2$    X :   $f(\mathbf{x}^c, \mathbf{w}), \ \forall \ c \in \mathbf{I}_1$     O :   $f(\mathbf{x}^c, \mathbf{w}), \ \forall \ c \in \mathbf{I}_2$



$$\alpha \equiv \min_{c \in \mathbf{I}_1} f(\mathbf{x}^c, \mathbf{w}) \ ; \ \beta \equiv \max_{c \in \mathbf{I}_2} f(\mathbf{x}^c, \mathbf{w})$$

34

# the classification inferencing mechanism



When LSC ($\alpha > \beta$) is satisfied, the classification inferencing mechanism

$$f(\mathbf{x}^c, \mathbf{w}) \geq v \ \forall \ c \in \mathbf{I}_1 \text{ and } f(\mathbf{x}^c, \mathbf{w}) \leq -v \ \forall \ c \in \mathbf{I}_2$$

can be set by directly adjusting $\mathbf{w}^o$ according to the following formula:

$$\frac{2v}{\alpha - \beta} w_i^o \rightarrow w_i^o \ \forall \ i,$$

The weight vector between the hidden layer and the output node

$$\text{then } v - \min_{c \in \mathbf{I}_1} \sum_{i=1}^p w_i^o a_i^c \rightarrow w_0^o$$
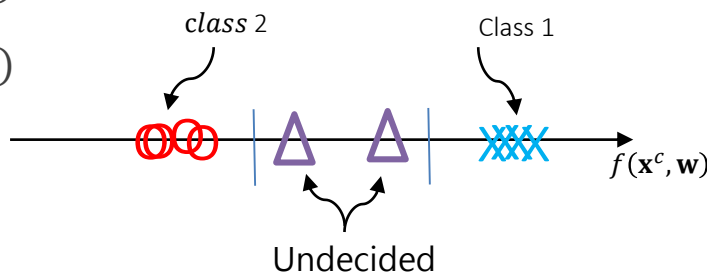
The threshold of the output node

35

# The classification applications

$y^c = \text{class } 1 \; \forall \; c \in \mathbf{I}_1(n); \; y^c = \text{class } 2 \; \forall \; c \in \mathbf{I}_2(n); \; \mathbf{I}(n) \equiv \mathbf{I}_1(n) \cup \mathbf{I}_2(n)$

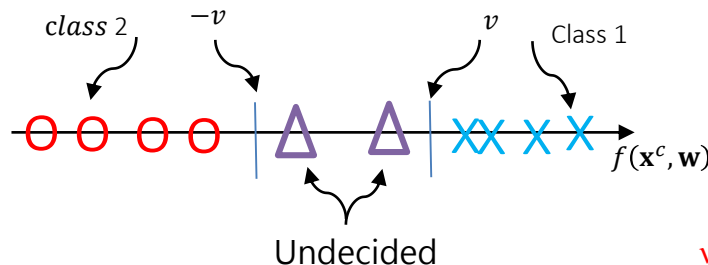**X**: $f(\mathbf{x}^c, \mathbf{w}) \; \forall \; c \in \mathbf{I}_1(n)$

**O**: $f(\mathbf{x}^c, \mathbf{w}) \; \forall \; c \in \mathbf{I}_2(n)$

*class* 2

Class 1

$f(\mathbf{x}^c, \mathbf{w})$

Undecided

learning goal type I:
$|f(\mathbf{x}^c, \mathbf{w}) - 1| \le \varepsilon \; \forall \; c \in \mathbf{I}_1(n);$
$|f(\mathbf{x}^c, \mathbf{w}) + 1| \le \varepsilon \; \forall \; c \in \mathbf{I}_2(n)$

ε Is a hyperparameter regarding the learning!

The inferencing mechanism

*class* 2

$-v$

$v$

Class 1

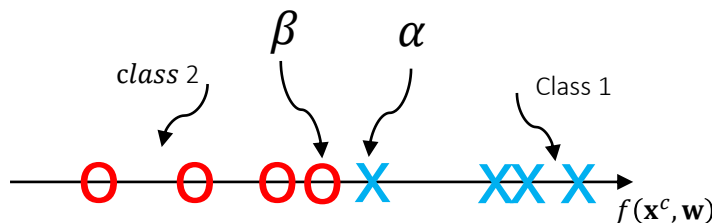$f(\mathbf{x}^c, \mathbf{w})$

Undecided

learning goal type II
(also inferencing goal):
$f(\mathbf{x}^c, \mathbf{w}) \ge v \; \forall \; c \in \mathbf{I}_1(n);$
$f(\mathbf{x}^c, \mathbf{w}) \le -v \; \forall \; c \in \mathbf{I}_2(n)$
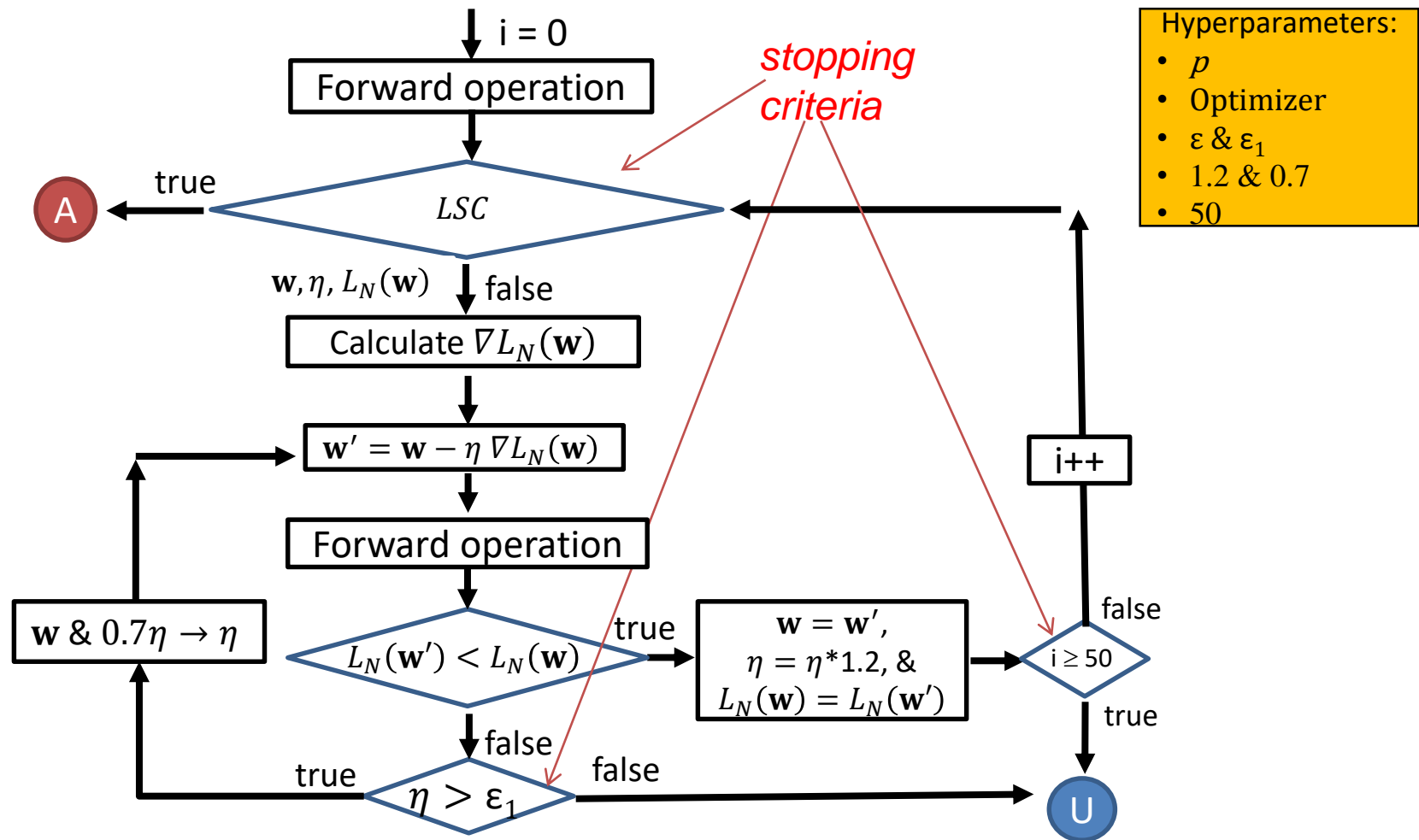
ν Is a hyperparameter regarding the inferencing!

$\alpha = \min\limits_{c \in \mathbf{I}_1(n)} f(\mathbf{x}^c, \mathbf{w})$

$\beta = \max\limits_{c \in \mathbf{I}_2(n)} f(\mathbf{x}^c, \mathbf{w})$

$\beta$

$\alpha$

*class* 2

Class 1

$f(\mathbf{x}^c, \mathbf{w})$

learning goal type III: LSC

# The flowchart of weight-tuning module_EU_LG_UA



*stopping criteria*

**Hyperparameters:**
- $p$
- Optimizer
- $\varepsilon$ & $\varepsilon_1$
- 1.2 & 0.7
- 50

i = 0

Forward operation

A — true

LSC

$\mathbf{w}, \eta, L_N(\mathbf{w})$   false

Calculate $\nabla L_N(\mathbf{w})$

$\mathbf{w}' = \mathbf{w} - \eta \, \nabla L_N(\mathbf{w})$

Forward operation

$\mathbf{w}$ & $0.7\eta \to \eta$

$L_N(\mathbf{w}') < L_N(\mathbf{w})$   true

$\mathbf{w} = \mathbf{w}'$, $\eta = \eta*1.2$, & $L_N(\mathbf{w}) = L_N(\mathbf{w}')$

false

$\eta > \varepsilon_1$   false

i++

$i \geq 50$   false
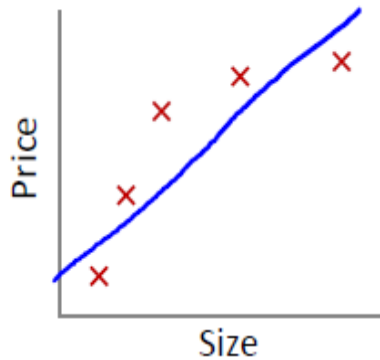
true

U

Where we are now...

# Algorithm representation and development

(Algorithm - Wikipedia)

- Algorithms can be expressed in many kinds of notation, including natural languages, pseudocode, flowcharts, drakon-charts, programming languages or control tables (processed by interpreters).

  ✓ Natural language expressions of algorithms tend to be verbose and ambiguous, and are rarely used for complex or technical algorithms.

  ✓ Pseudocode, flowcharts, drakon-charts and control tables are structured ways to express algorithms that avoid many of the ambiguities common in the statements based on natural language.

  ✓ Programming languages are primarily intended for expressing algorithms in a form that can be executed by a computer, but are also often used as a way to define or document algorithms.

- Typical steps in the development of algorithms:
  ✓ Problem definition     ← learning-based prediction problem
  ✓ Development of a model   ← 2-layer net, 4-layer net, or deep neural networks
  ✓ Specification of the algorithm   ← The learning algorithm
  ✓ Designing an algorithm   ← The gradient-descent-based learning algorithm
  ✓ Checking the correctness of the algorithm   ← The math proof of the proposed learning algorithm
  ✓ Analysis of algorithm   ← The amount of parameters, the (learning and inferencing) time scale, …
  ✓ Implementation of algorithm  ← The coding
  ✓ Program testing
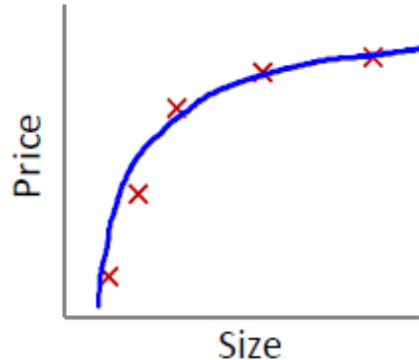  ✓ Documentation preparation

38

# Program testing --
# Performance of AI Applications

- How do AI professionals evaluate the performance of the AI applications?

  ← effectiveness & efficiency

- However, there are learning dilemma and overfitting when evaluating the effectiveness & efficiency.

- You need to deal with learning dilemma and overfitting, not only for the purposes of learning, but also of inferencing.
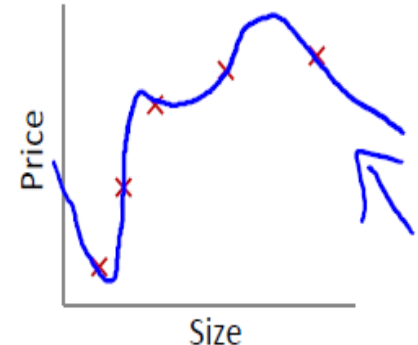
# overfitting



$\rightarrow \theta_0 + \theta_1 x$

"Undefit"   "High bias"

$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$

"Just right"

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

"Overfit"   "High variance"

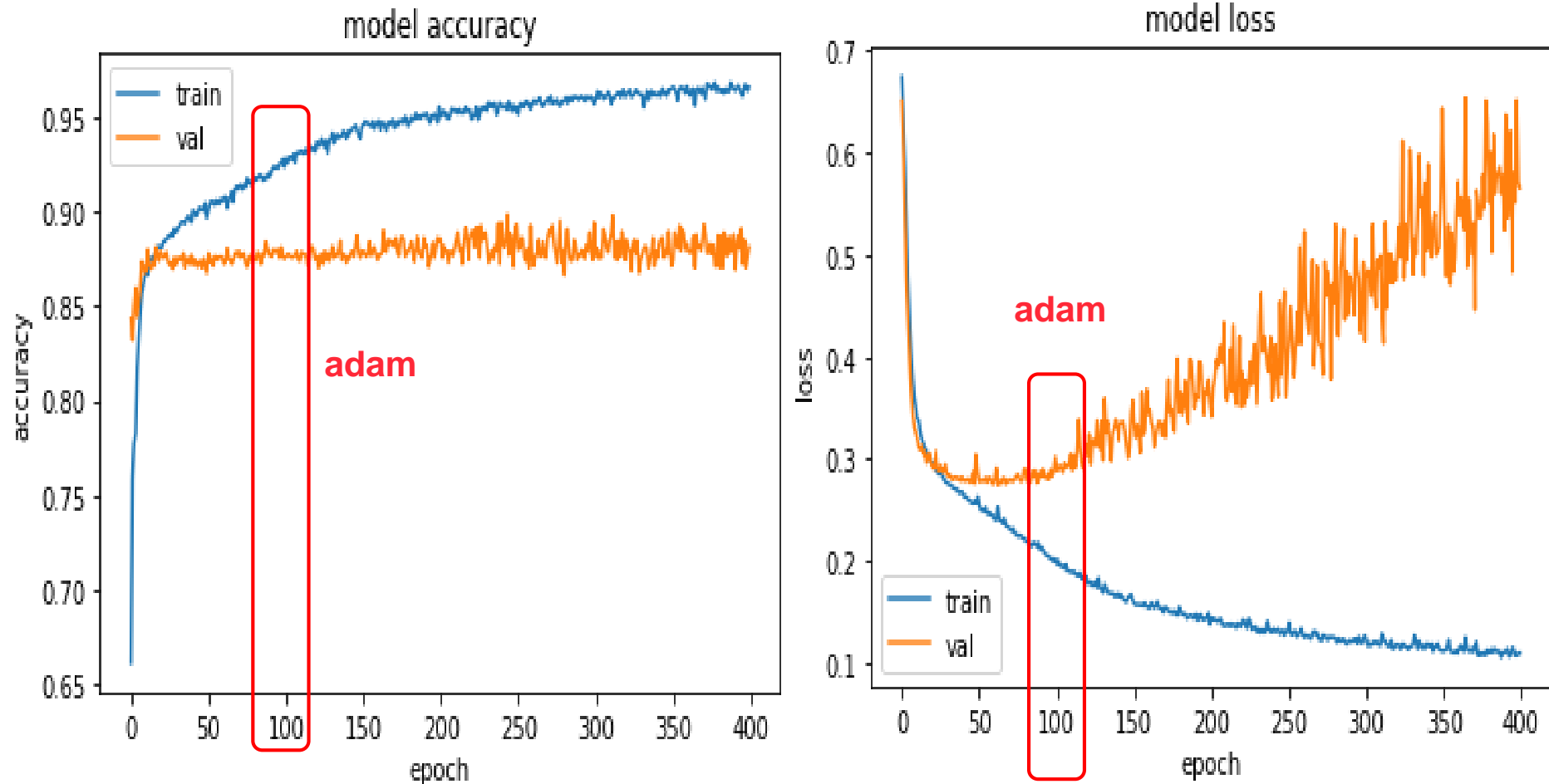inadequate          good compromise          over-fitting

# Generalization

- Learned hypothesis may fit the training data very well, even noises (or outliers) in the training data, but fail to generalize to new examples (test data)
- In machine learning and statistical learning, the generalization error (also known as the out-of-sample error) is a measure of how accurately an algorithm is able to predict outcome values for previously unseen data.

# Learning curves

- Because learning algorithms are evaluated on finite samples, the evaluation of a learning algorithm may be sensitive to sampling error.
- As a result, measurements of prediction error on the current data may not provide much information about predictive ability on new data.
- The performance of a learning algorithm is measured by plots of the generalization error values through the learning process, which are called learning curves.
- Generalization error can be minimized by avoiding overfitting in the learning process.

# Learning curve and overfitting
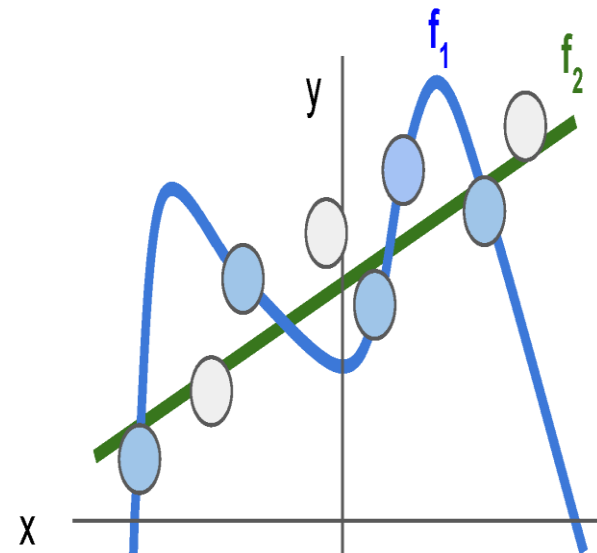


model accuracy

model loss

adam

adam

Early stop!

# Overfitting

In statistics, **overfitting** is "the production of an analysis that corresponds too closely or exactly to a particular set of data, and may therefore fail to fit additional data or predict future observations reliably."
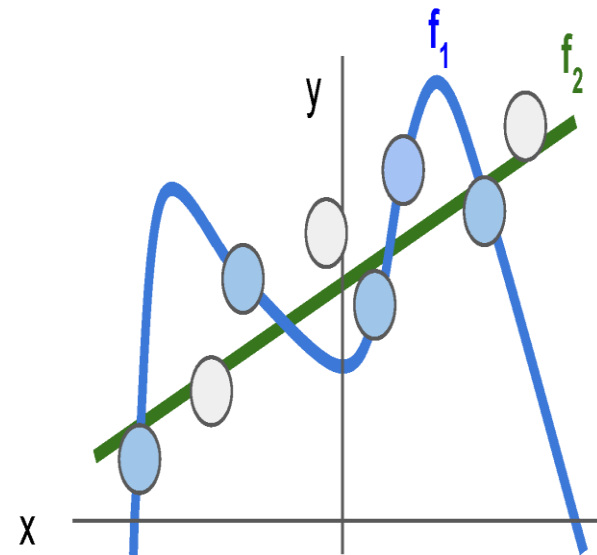
## Regularization: Prefer Simpler Models



Regularization pushes against fitting the data *too* well so we don't fit noise in the data

# Overfitting

An **over-fitted model** is a model that contains <span style="color:red">more parameters</span> than can be justified by the data.

Regularization: Prefer Simpler Models



Regularization pushes against fitting the data *too* well so we don't fit noise in the data

# Regularization

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W)$$

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

**Occam's Razar**: Among multiple competing hypotheses, the simplest is the best, William of Ockham 1285-1347

# Regularization

$\lambda$ = regularization strength (hyperparameter)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W)$$

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

**Simple examples**

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

**More complex**:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

# Regularization

$\lambda$ = regularization strength (hyperparameter)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W)$$

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

Why regularize?
- Express preferences over weights
- Make the model *simple* so it works on test data
- Improve optimization by adding curvature

# Regularization - In practice

**Training**: Add random noise
**Testing**: Marginalize over the noise

**Examples**:
Dropout
Batch Normalization
Data Augmentation
DropConnect
Fractional Max Pooling
Stochastic Depth
Cutout / Random Crop
Mixup

- Consider dropout for large fully-connected layers
- Batch normalization and data augmentation almost always a good idea
- Try cutout and mixup especially for small classification datasets
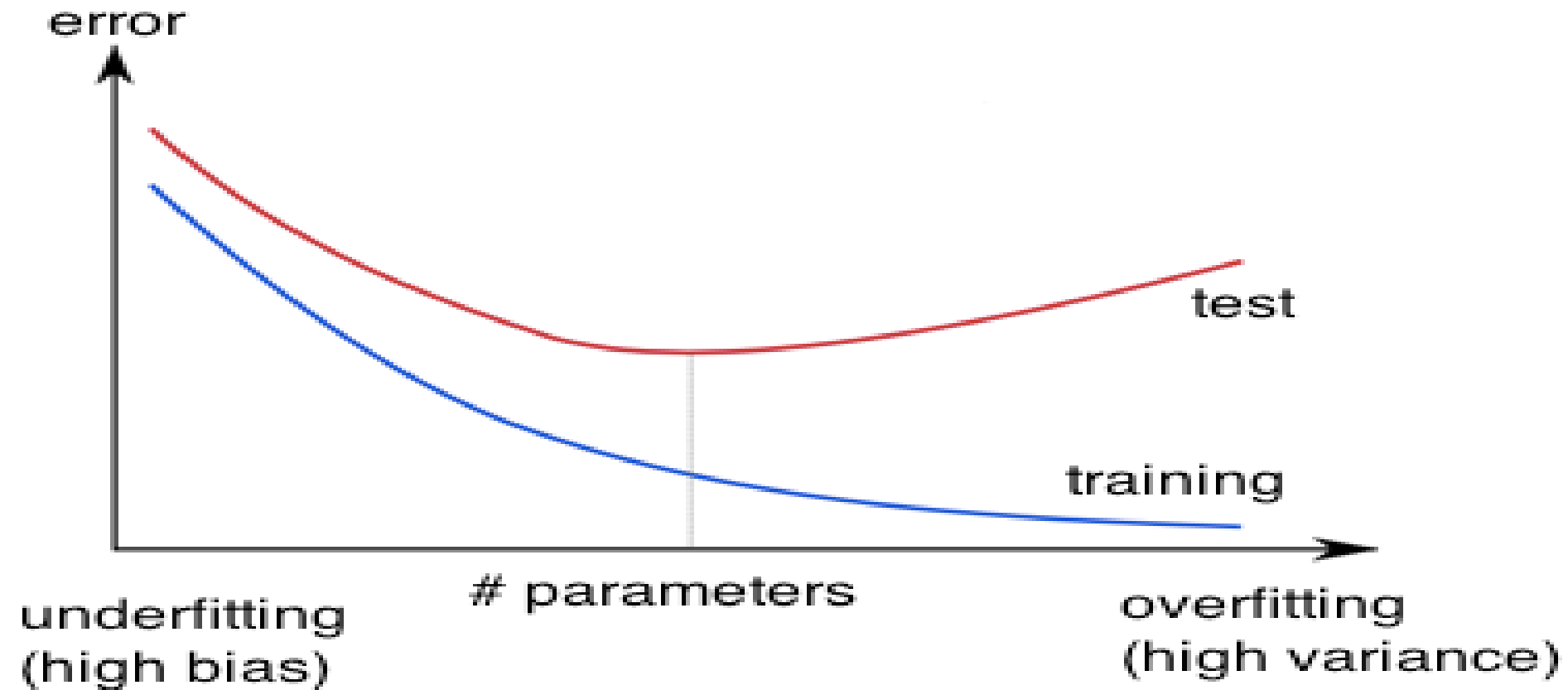
# Summary: the overfitting may be due to big weights

- Adopt a regularization term in the loss function to penalize big weights:

$$\mathrm{L}_N(\mathbf{w}) \equiv \frac{1}{N}\sum_{c=1}^{N}(f(\mathbf{x}^c, \mathbf{w}) - y^c)^2 + \lambda\|\mathbf{w}\|^2$$
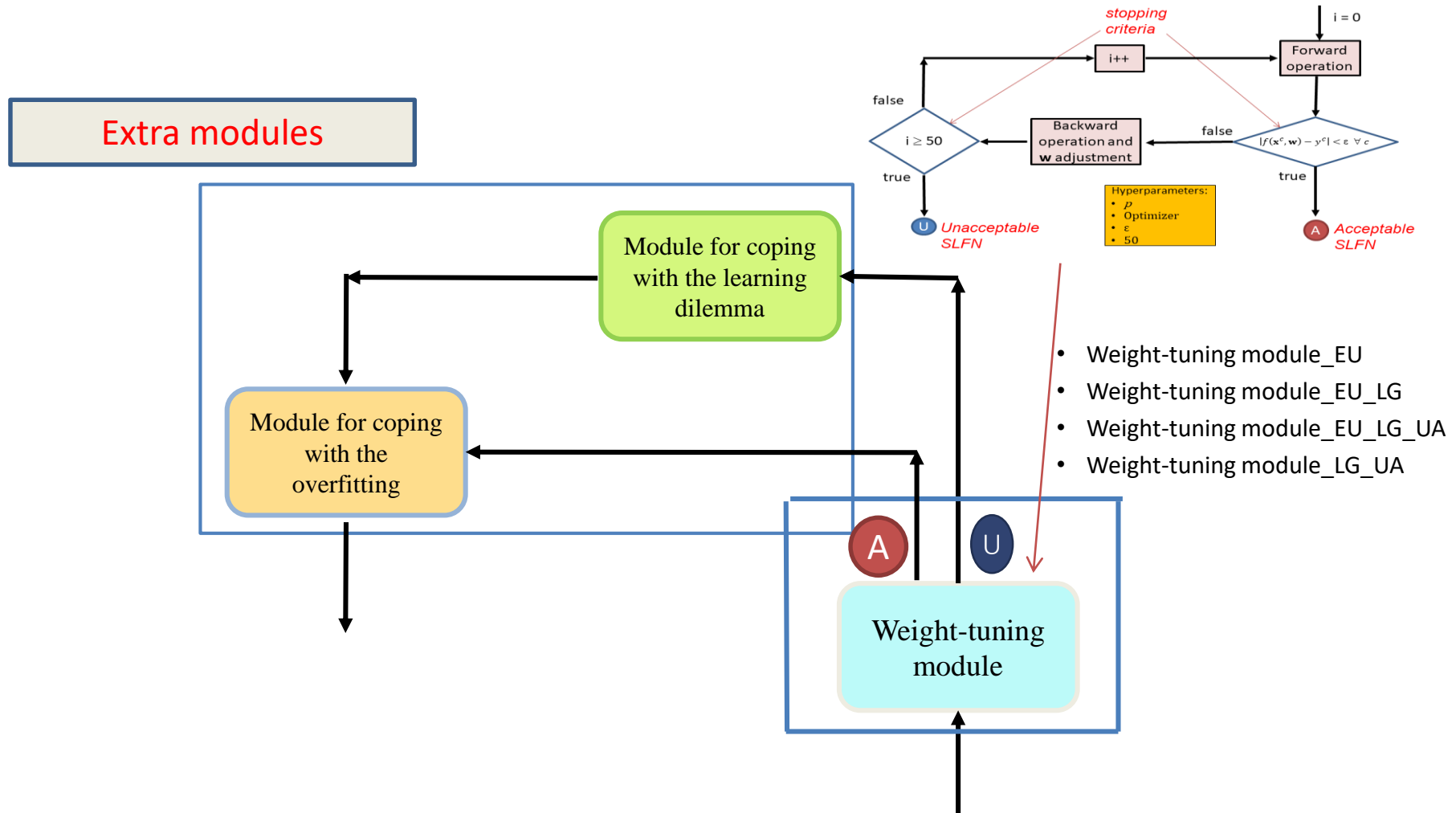
- The regularization coefficient $\lambda$ determines how dominant the regularization is during gradient computation

- Big regularization coefficient → big penalty for big weights

- The above is the L2 regularization

- L1 regularization: $\lambda|\mathbf{w}|$

- Elastic net regularization: L1 + L2

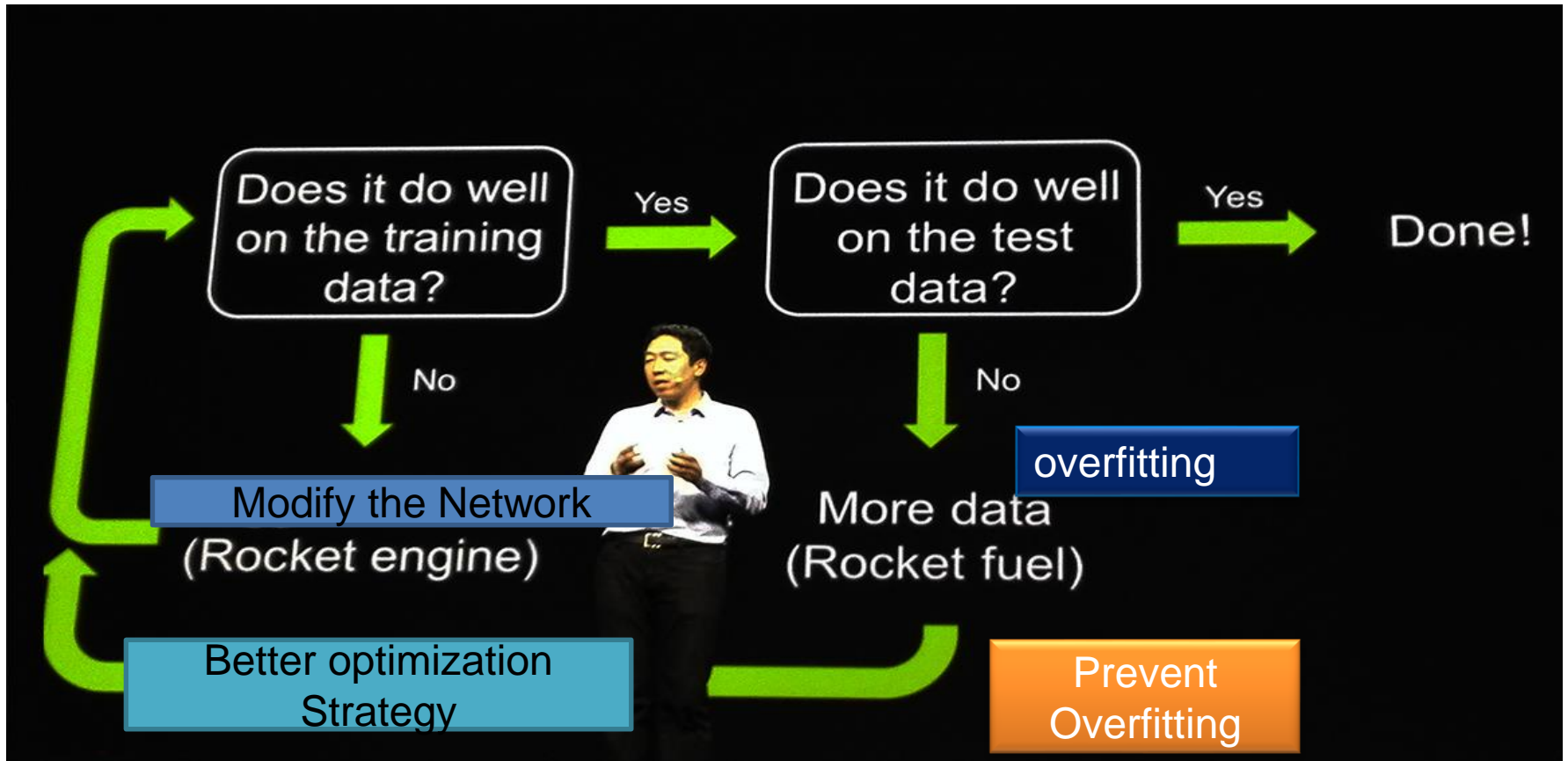# Summary: the overfitting may be due to too many hidden nodes



*https://www.neuraldesigner.com/images/learning/selection_error.svg*

# Inferencing Issues

Extra modules

Module for coping with the learning dilemma

Module for coping with the overfitting

Weight-tuning module

A

U

stopping criteria

i = 0

i++

Forward operation

false

false

$i \geq 50$

Backward operation and **w** adjustment

$|f(\mathbf{x}^c, \mathbf{w}) - y^c| < \varepsilon \ \forall \ c$

true

true

U *Unacceptable SLFN*

A *Acceptable SLFN*

Hyperparameters:
- $p$
- Optimizer
- $\varepsilon$
- 50

- Weight-tuning module_EU
- Weight-tuning module_EU_LG
- Weight-tuning module_EU_LG_UA
- Weight-tuning module_LG_UA

# Recipe for Deep Learning



http://www.gizmodo.com.au/2015/04/the-basic-recipe-for-machine-learning-explained-in-a-single-powerpoint-slide/

# Dealing with the overfitting due to big weights – the regularizing module

$$L_N(\mathbf{w}) \equiv \frac{\sum_c (f(\mathbf{x}^c, \mathbf{w}) - y^c)^2}{N} + \lambda \left( \sum_{i=0}^{p} (w_i^o)^2 + \sum_{i=1}^{p} \sum_{j=0}^{m} (w_{ij}^H)^2 \right)$$



Module for coping with the learning dilemma

Module for coping with the overfitting

A  U

Weight-tuning module