# Machine Learning Foundations (機器學習基石)

Lecture 13: Hazard of Overfitting

Hsuan-Tien Lin (林軒田) htlin@csie.ntu.edu.tw

Department of Computer Science & Information Engineering

National Taiwan University (國立台灣大學資訊工程系)



# Roadmap

- 1 When Can Machines Learn?
- 2 Why Can Machines Learn?
- 3 How Can Machines Learn?

Lecture 12: Nonlinear Transform

nonlinear uia nonlinear feature transform Φ plus linear with price of model complexity

4 How Can Machines Learn Better?

#### Lecture 13: Hazard of Overfitting

- What is Overfitting?
- The Role of Noise and Data Size
- Deterministic Noise
- Dealing with Overfitting

What is Overfitting?

#### **Bad Generalization**

- regression for *x* ∈ ℝ with *N* = 5 examples
- target f(x) = 2nd order polynomial
- label  $y_n = f(x_n) + \text{very small noise}$
- linear regression in *Z*-space +
  Φ = 4th order polynomial
- unique solution passing all examples  $\implies E_{in}(g) = 0$
- *E*<sub>out</sub>(*g*) huge

bad generalization: low *E*<sub>in</sub>, high *E*<sub>out</sub>



#### Hazard of Overfitting

What is Overfitting?

# Bad Generalization and Overfitting

- take d<sub>VC</sub> = 1126 for learning: bad generalization --(E<sub>out</sub> - E<sub>in</sub>) large
- switch from d<sub>vc</sub> = d<sup>\*</sup><sub>vc</sub> to d<sub>vc</sub> = 1126:
  overfitting
  -E<sub>in</sub> ↓, E<sub>out</sub> ↑
- switch from d<sub>vc</sub> = d<sup>\*</sup><sub>vc</sub> to d<sub>vc</sub> = 1: underfitting —E<sub>in</sub> ↑, E<sub>out</sub> ↑



bad generalization: low *E*<sub>in</sub>, high *E*<sub>out</sub>; overfitting: lower *E*<sub>in</sub>, higher *E*<sub>out</sub>



# next: how does **noise** & **data size** affect overfitting?

Hsuan-Tien Lin (NTU CSIE)

Based on our discussion, for data of fixed size, which of the following situation is relatively of the lowest risk of overfitting?

- **1** small noise, fitting from small  $d_{\rm vc}$  to median  $d_{\rm vc}$
- small noise, fitting from small  $d_{\rm VC}$  to large  $d_{\rm VC}$ 2
- 3 large noise, fitting from small  $d_{\rm VC}$  to median  $d_{\rm VC}$
- 4 large noise, fitting from small  $d_{\rm VC}$  to large  $d_{\rm VC}$

Based on our discussion, for data of fixed size, which of the following situation is relatively of the lowest risk of overfitting?

- **1** small noise, fitting from small  $d_{vc}$  to median  $d_{vc}$
- 2 small noise, fitting from small d<sub>VC</sub> to large d<sub>VC</sub>
- **(3)** large noise, fitting from small  $d_{VC}$  to median  $d_{VC}$
- Iarge noise, fitting from small d<sub>VC</sub> to large d<sub>VC</sub>

#### Reference Answer: (1)

Two causes of overfitting are noise and excessive  $d_{VC}$ . So if both are relatively 'under control', the risk of overfitting is smaller.

The Role of Noise and Data Size

#### Case Study (1/2)



#### overfitting from best $g_2 \in \mathcal{H}_2$ to best $g_{10} \in \mathcal{H}_{10}$ ?

Hsuan-Tien Lin (NTU CSIE)

The Role of Noise and Data Size

#### Case Study (2/2)



overfitting from  $g_2$  to  $g_{10}$ ? both yes!

Hsuan-Tien Lin (NTU CSIE)

Hazard of Overfitting

#### The Role of Noise and Data Size Irony of Two Learners



#### but *R* wins in *E*<sub>out</sub> a lot! philosophy: concession for advantage? :-)

Hsuan-Tien Lin (NTU CSIE)

Hazard of Overfitting

The Role of Noise and Data Size

## Learning Curves Revisited



- $\mathcal{H}_{10}$ : lower  $\overline{E_{out}}$  when  $N \to \infty$ , but much larger generalization error for small N
- gray area : *O* overfits!  $(\overline{E_{in}} \downarrow, \overline{E_{out}} \uparrow)$

#### *R* always wins in $\overline{E_{out}}$ if *N* small!

Hsuan-Tien Lin (NTU CSIE)

#### Hazard of Overfitting

The Role of Noise and Data Size

## The 'No Noise' Case



Hsuan-Tien Lin (NTU CSIE)

When having limited data, in which of the following case would learner R perform better than learner O?

- 1 limited data from a 10-th order target function with some noise
- 2 limited data from a 1126-th order target function with no noise
- 3 limited data from a 1126-th order target function with some noise
- 4 all of the above

When having limited data, in which of the following case would learner R perform better than learner O?

- 1 limited data from a 10-th order target function with some noise
- 2 limited data from a 1126-th order target function with no noise
- 3 limited data from a 1126-th order target function with some noise
- 4 all of the above

#### Reference Answer: (4)

We discussed about (1) and (2), but you shall be able to 'generalize' :-) that *R* also wins in the more difficult case of (3).

Deterministic Noise

#### A Detailed Experiment





- Gaussian iid noise  $\epsilon$  with level  $\sigma^2$
- some 'uniform' distribution on f(x) with complexity level Q<sub>f</sub>

• data size N

goal: **'overfit level'** for different  $(N, \sigma^2)$  and  $(N, Q_f)$ ?

Deterministic Noise

#### The Overfit Measure



0

overfit measure  $E_{out}(g_{10}) - E_{out}(g_2)$ 

#### Hazard of Overfitting

Deterministic Noise

## The Results



Hsuan-Tien Lin (NTU CSIE)



# Impact of Noise and Data Size



#### overfitting 'easily' happens

Hsuan-Tien Lin (NTU CSIE)

#### **Deterministic Noise**

- if *f* ∉ *H*: something of *f* cannot be captured by *H*
- deterministic noise : difference between best  $h^* \in \mathcal{H}$  and f
- acts like 'stochastic noise'—not new to CS: pseudo-random generator
- difference to stochastic noise:
  - depends on *H*
  - fixed for a given x



philosophy: when teaching a kid, perhaps better not to use examples from a complicated target function? :-)

Consider the target function being sin(1126x) for  $x \in [0, 2\pi]$ . When x is uniformly sampled from the range, and we use all possible linear hypotheses  $h(x) = w \cdot x$  to approximate the target function with respect to the squared error, what is the level of deterministic noise for each x?

1 | sin(1126*x*)|

**2** 
$$|\sin(1126x) - x|$$

- **3**  $|\sin(1126x) + x|$
- $4 |\sin(1126x) 1126x|$

Consider the target function being sin(1126x) for  $x \in [0, 2\pi]$ . When x is uniformly sampled from the range, and we use all possible linear hypotheses  $h(x) = w \cdot x$  to approximate the target function with respect to the squared error, what is the level of deterministic noise for each x?

1 | sin(1126*x*)|

**2** 
$$|\sin(1126x) - x|$$

- **3**  $|\sin(1126x) + x|$
- $4 |\sin(1126x) 1126x|$

#### Reference Answer: (1)

You can try a few different *w* and convince yourself that the best hypothesis  $h^*$  is  $h^*(x) = 0$ . The deterministic noise is the difference between *f* and  $h^*$ . Dealing with Overfitting

## Driving Analogy Revisited

learning	driving
overfit	commit a car accident
use excessive $d_{VC}$	'drive too fast'
noise	bumpy road
limited data size N	limited observations about road condition
start from simple model	drive slowly
data cleaning/pruning	use more accurate road information
data hinting	exploit more road information
regularization	put the brakes

#### all very practical techniques to combat overfitting

#### Hazard of Overfitting

Dealing with Overfitting

# Data Cleaning/Pruning



- if 'detect' the outlier 5 at the top by

  - wrong by current classifier
  - . . .
- possible action 1: correct the label (data cleaning)
- possible action 2: remove the example (data pruning)

#### possibly helps, but effect varies

Hsuan-Tien Lin (NTU CSIE)

# Data Hinting



- slightly shifted/rotated digits carry the same meaning
- possible action: add virtual examples by shifting/rotating the given digits (data hinting)

Hsuan-Tien Lin (NTU CSIE)

Assume we know that f(x) is symmetric for some 1D regression application. That is, f(x) = f(-x). One possibility of using the knowledge is to consider symmetric hypotheses only. On the other hand, you can also generate virtual examples from the original data  $\{(x_n, y_n)\}$  as hints. What virtual examples suit your needs best?

$$\{(x_n,-y_n)\}$$

**2** 
$$\{(-x_n, -y_n)\}$$

**3** 
$$\{(-x_n, y_n)\}$$

$$4 \{(2x_n, 2y_n)\}$$

Assume we know that f(x) is symmetric for some 1D regression application. That is, f(x) = f(-x). One possibility of using the knowledge is to consider symmetric hypotheses only. On the other hand, you can also generate virtual examples from the original data  $\{(x_n, y_n)\}$  as hints. What virtual examples suit your needs best?

1 
$$\{(x_n, -y_n)\}$$

**2** 
$$\{(-x_n, -y_n)\}$$

**3** 
$$\{(-x_n, y_n)\}$$

$$4 \{(2x_n, 2y_n)\}$$

#### Reference Answer: (3)

We want the virtual examples to encode the invariance when  $x \rightarrow -x$ .



Dealing with Overfitting

- When Can Machines Learn?
- 2 Why Can Machines Learn?
- 3 How Can Machines Learn?

Lecture 12: Nonlinear Transform

4 How Can Machines Learn Better?

Lecture 13: Hazard of Overfitting

• What is Overfitting?

lower Ein but higher Eout

- The Role of Noise and Data Size overfitting 'easily' happens!
- Deterministic Noise what H cannot capture acts like noise
- Dealing with Overfitting

data cleaning/pruning/hinting, and more

next: putting the brakes with regularization