

# Machine Learning

## (機器學習)

### Lecture 07: Combatting 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

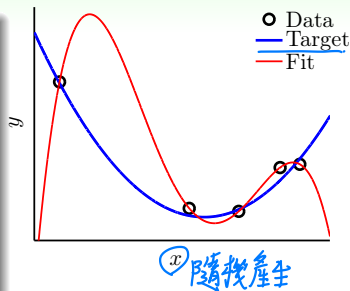
- 1 When Can Machines Learn?
- 2 Why Can Machines Learn?
- 3 How Can Machines Learn?
- 4 How Can Machines Learn **Better**?

## Lecture 07: Combatting Overfitting

- What is Overfitting?
- The Role of Noise and Data Size
- Deterministic Noise
- Dealing with Overfitting
- Regularized Hypothesis Set
- Weight Decay Regularization
- Regularization and VC Theory
- General Regularizers

# Bad Generalization

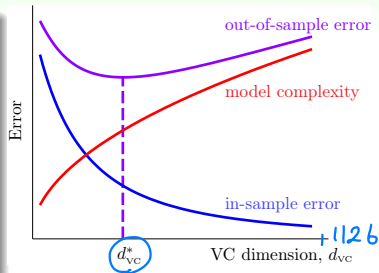
- regression for  $x \in \mathbb{R}$  with  $N = 5$  examples
- target  $f(x) = 2\text{nd order polynomial}$
- label  $y_n = f(x_n) + \text{very small noise}$
- linear regression in  $\mathcal{Z}$ -space +  $\Phi = 4\text{th order polynomial}$
- unique solution passing all examples  
 $\implies \underline{E_{\text{in}}(g)} = 0$
- $E_{\text{out}}(g)$  huge  $E_{\text{in}}, E_{\text{out}}$  差很多



bad generalization: low  $E_{\text{in}}$ , high  $E_{\text{out}}$

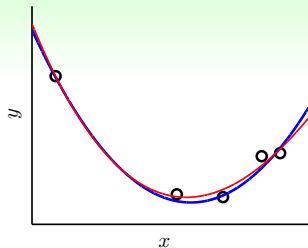
# Bad Generalization and Overfitting

- take  $d_{VC} = 1126$  for learning:  
bad generalization  
—( $E_{out} - E_{in}$ ) large 很大的數字.
- switch from  $d_{VC} = d_{VC}^*$  to  $d_{VC} = 1126$ :  
**overfitting**  
— $E_{in} \downarrow$ ,  $E_{out} \uparrow$  ← overfitting 發生
- switch from  $d_{VC} = d_{VC}^*$  to  $d_{VC} = 1$ :  
**underfitting** ← 做不夠好.  
— $E_{in} \uparrow$ ,  $E_{out} \uparrow$

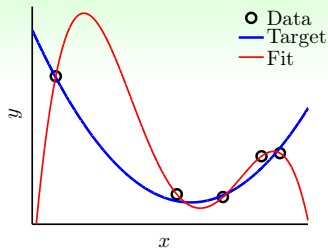


bad generalization: low  $E_{in}$ , high  $E_{out}$ ;  
**overfitting**: lower  $E_{in}$ , higher  $E_{out}$

# Cause of Overfitting: A Driving Analogy



'good fit'



**overfit**

learning

driving

overfit *mode complexity 太大*

commit a car accident

use excessive  $d_{VC}$

'drive too fast'

noise

*資料量太小*

bumpy road

limited data size  $N$

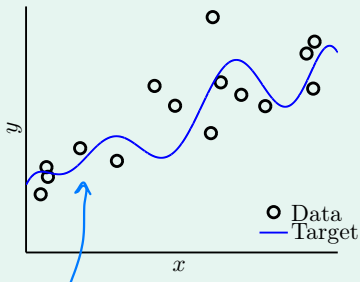
limited observations about road condition

next: how does noise & data size affect overfitting?

**Questions?**

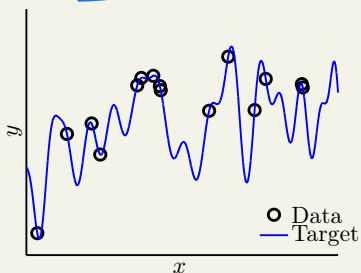
## Case Study (1/2)

10-th order target function  
+ noise



data 有 noise

50-th order target function  
noiselessly



沒有 noise.

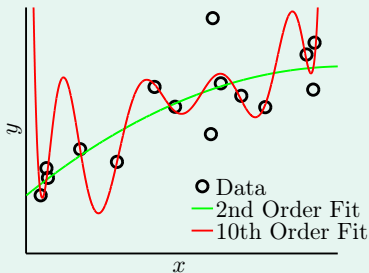
2次多項式  $\in \mathcal{H}_2$

10次多項式  $\in \mathcal{H}_{10}$

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

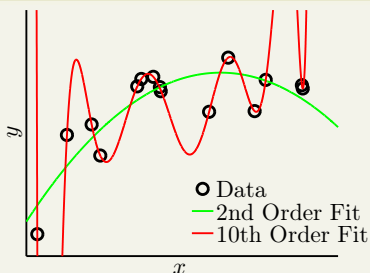
# Case Study (2/2)

## 10-th order target function + noise



	$g_2 \in \mathcal{H}_2$	$g_{10} \in \mathcal{H}_{10}$
$E_{in}$	0.050	0.034 较小
$E_{out}$	<u>0.127</u>	<u>9.00</u> over fit!

## 50-th order target function noiselessly

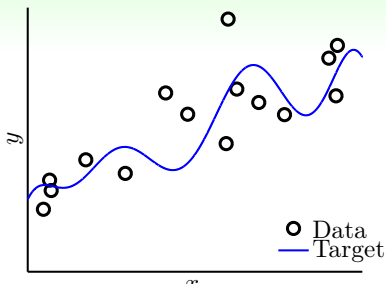


	$g_2 \in \mathcal{H}_2$	$g_{10} \in \mathcal{H}_{10}$
$E_{in}$	0.029	<u>0.00001</u>
$E_{out}$	0.120	<u>7680</u> overfit

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



# Irony of Two Learners

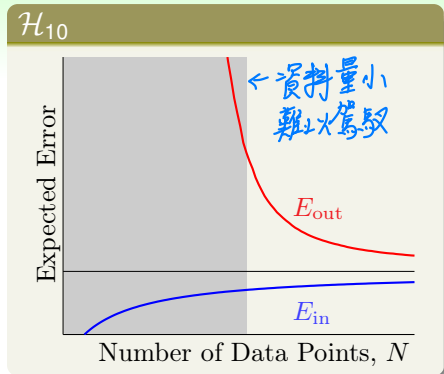
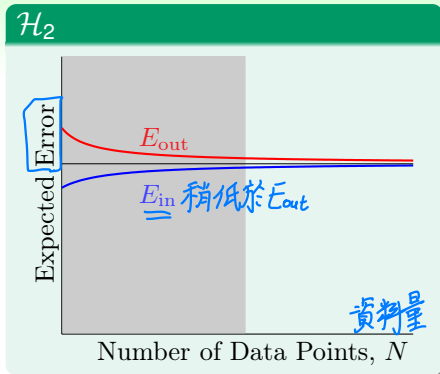


- learner **Overfit**: pick  $g_{10} \in \mathcal{H}_{10}$
- learner **Restrict**: pick  $g_2 \in \mathcal{H}_2$
- when both **know that target = 10th**  
—  $R$  'gives up' ability to fit

弱學生做的比較好。

but  $R$  **wins in  $E_{out}$**  a lot! 小區為進  
philosophy: concession for **advantage**? :-)

# Learning Curves Revisited

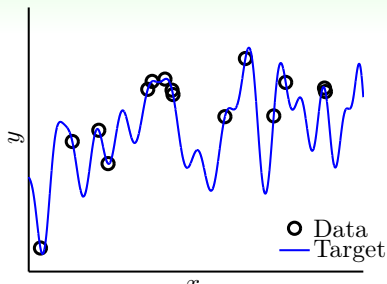
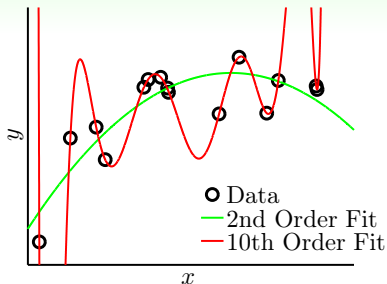


- $\mathcal{H}_{10}$ : lower  $\overline{E_{out}}$  when  $N \rightarrow \infty$ , but much larger generalization error for small  $N$
- gray area:  $O$  overfits! ( $\overline{E_{in}} \downarrow$ ,  $\overline{E_{out}} \uparrow$ )

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

# The 'No Noise' Case

沒有 noise  $\mathcal{H}_2$  還是做不好?



- learner **Overfit**: pick  $g_{10} \in \mathcal{H}_{10}$
- learner **Restrict**: pick  $g_2 \in \mathcal{H}_2$
- when both **know that there is no noise** —  $R$  still wins

is there really **no noise**?  
'target complexity' acts like noise

Target 真的太複雜

# Questions?

# A Detailed Experiment

$$y = f(x) + \epsilon$$

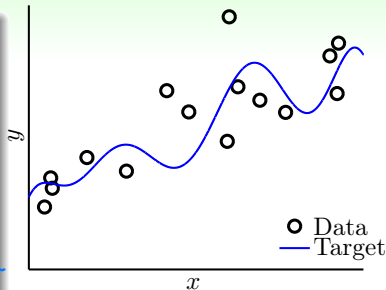
gaussian noise.

$$\sim \text{Gaussian} \left( \underbrace{\sum_{q=0}^{Q_f} \alpha_q x^q}_{f(x)}, \sigma^2 \right)$$

noise level

- Gaussian iid noise  $\epsilon$  with level  $\sigma^2$
- some 'uniform' distribution on  $f(x)$  with complexity level  $Q_f$  某次方的 target function
- data size  $N$

$\sigma^2$  与  $Q_f$  对 overfit 的影响

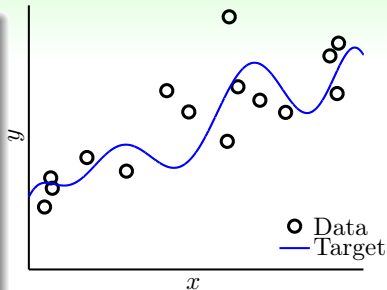


goal: 'overfit level' for different  $(N, \sigma^2)$  and  $(N, Q_f)$ ?

# The Overfit Measure



- $g_2 \in \mathcal{H}_2$
- $g_{10} \in \mathcal{H}_{10}$
- $E_{in}(g_{10}) \leq E_{in}(g_2)$  for sure



量測 overfit

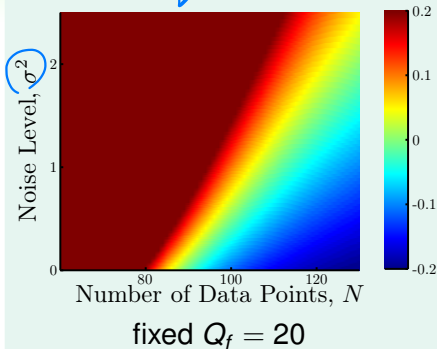


overfit measure  $E_{out}(g_{10}) - E_{out}(g_2)$

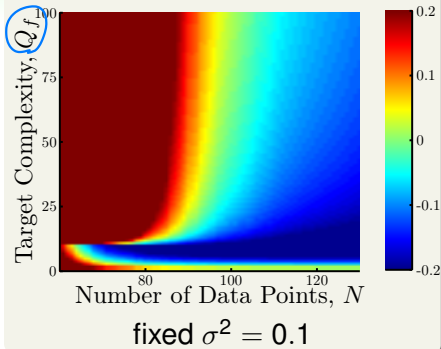
紅色 overfit

# The Results

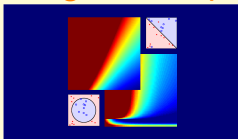
impact of  $\sigma^2$  versus  $N$



impact of  $Q_f$  versus  $N$



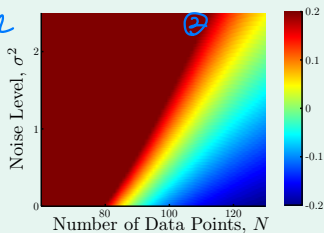
ring a bell? :-)



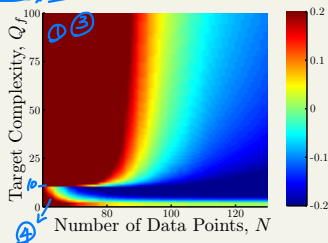
# Impact of Noise and Data Size

impact of  $\sigma^2$  versus  $N$ :  
stochastic noise

随机



impact of  $Q_f$  versus  $N$ :  
deterministic noise



four reasons of serious overfitting:

① data size $N \downarrow$	overfit $\uparrow$
② stochastic noise $\uparrow$	overfit $\uparrow$
③ deterministic noise $\uparrow$	overfit $\uparrow$
④ excessive power $\uparrow$	overfit $\uparrow$

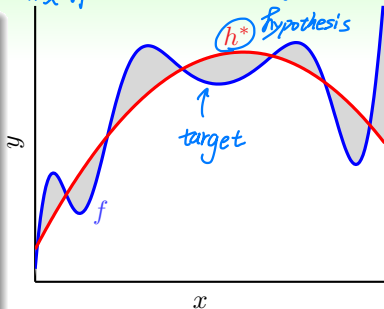
overfitting 'easily' happens



# Deterministic Noise

\* target function 太複雜 跟 noise 沒啥 2 樣.

- if  $f \notin \mathcal{H}$ : something of  $f$  cannot be captured by  $\mathcal{H}$
- **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  $\mathcal{H}$
  - fixed for a given  $x$



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

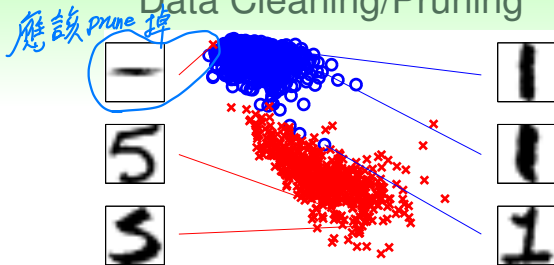
**Questions?**

# 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
★ <b>start from simple model</b> ★ <b>data cleaning/pruning</b> <small>整理数据</small> <b>data hinting</b> ✓ <b>regularization</b> <u><b>validation</b></u>	<u>drive slowly</u> use more accurate road information exploit more road information put the brakes monitor the dashboard

all very **practical** techniques  
to combat overfitting

# Data Cleaning/Pruning



- if 'detect' the outlier **5** at the top by
  - too close to other  $\circ$ , or too far from other  $\times$
  - wrong by current classifier
  - ... *eg. daytime  $\rightarrow$  nighttime*
- possible action 1: correct the label (data cleaning)
- possible action 2: remove the example (data pruning)

possibly helps, but effect varies

# Data Hinting *data augmentation*



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

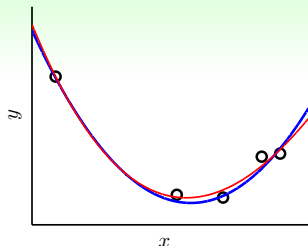
*不是同个 distribution!?*

possibly helps, but **watch out**

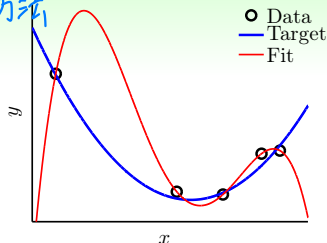
—**virtual example** not  $\overset{iid}{\sim} P(x, y)$ !

# Regularization: The Magic of 'Brake'

↑ 对付overfitting的方法

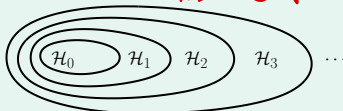


'regularized fit'



overfit

- idea: 'step back' from  $\mathcal{H}_{10}$  to  $\mathcal{H}_2$  高次走向低次. regularization

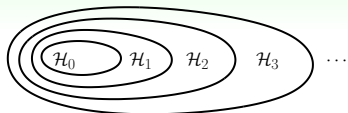


- name history: function approximation for **ill-posed problems**

how to step back?

# Questions?

# Stepping Back as Constraint



Q-th order polynomial transform for  $x \in \mathbb{R}$ :

$$\Phi_Q(x) = (1, x, x^2, \dots, x^Q)$$

+ linear regression, denote  $\tilde{\mathbf{w}}$  by  $\mathbf{w}$

hypothesis  $\mathbf{w}$  in  $\mathcal{H}_{10}$ :  $w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_{10} x^{10}$

hypothesis  $\mathbf{w}$  in  $\mathcal{H}_2$ :  $w_0 + w_1 x + w_2 x^2$

that is,  $\mathcal{H}_2 = \mathcal{H}_{10}$  AND constraint that  $w_3 = w_4 = \dots = w_{10} = 0$

step back = constraint



# Regression with Constraint

$$\mathcal{H}_{10} \equiv \left\{ \mathbf{w} \in \mathbb{R}^{10+1} \right\}$$

regression with  $\mathcal{H}_{10}$ :

$$\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w})$$

$$\mathcal{H}_2 \equiv \left\{ \mathbf{w} \in \mathbb{R}^{10+1} \right. \\ \left. \text{while } \underline{w_3 = w_4 = \dots = w_{10} = 0} \right\}$$

*constraint*

regression with  $\mathcal{H}_2$ :

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^{10+1}} \quad & E_{\text{in}}(\mathbf{w}) \\ \text{s.t.} \quad & \underline{w_3 = w_4 = \dots = w_{10} = 0} \end{aligned}$$

step back = **constrained optimization** of  $E_{\text{in}}$

why don't you just use  $\mathbf{w} \in \mathbb{R}^{2+1}$ ? :-)

# Regression with Looser Constraint

$$\mathcal{H}_2 \equiv \left\{ \mathbf{w} \in \mathbb{R}^{10+1} \right. \\ \left. \text{while } w_3 = \dots = w_{10} = 0 \right\}$$

regression with  $\mathcal{H}_2$ :

$$\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w})$$

$$\text{s.t. } w_3 = \dots = w_{10} = 0$$

放宽

$$\mathcal{H}'_2 \equiv \left\{ \mathbf{w} \in \mathbb{R}^{10+1} \right. \\ \left. \text{while } \geq 8 \text{ of } w_q = 0 \right\}$$

只要有8个w是0即可

regression with  $\mathcal{H}'_2$ :

$$\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w})$$

$$\text{s.t. } \sum_{q=0}^{10} \mathbb{I}[w_q \neq 0] \leq 3$$

- more flexible than  $\mathcal{H}_2$ :

$$\mathcal{H}_2 \subset \mathcal{H}'_2$$

- less risky than  $\mathcal{H}_{10}$ :

$$\mathcal{H}'_2 \subset \mathcal{H}_{10}$$

有3个系数可非零

bad news for sparse hypothesis set  $\mathcal{H}'_2$ :

**NP-hard to solve :-)**

只有少数系数非零

# Regression with Softer Constraint

還是不好解

$$\mathcal{H}'_2 \equiv \left\{ \mathbf{w} \in \mathbb{R}^{10+1} \right. \\ \left. \text{while } \geq 8 \text{ of } w_q = 0 \right\}$$

regression with  $\mathcal{H}'_2$ :

$$\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w}) \text{ s.t. } \sum_{q=0}^{10} \mathbb{I}[w_q \neq 0] \leq 3$$

$$\mathcal{H}(C) \equiv \left\{ \mathbf{w} \in \mathbb{R}^{10+1} \right. \\ \left. \text{while } \|\mathbf{w}\|^2 \leq C \right\}$$

regression with  $\mathcal{H}(C)$ :

$$\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w}) \text{ s.t. } \sum_{q=0}^{10} w_q^2 \leq C$$

(w的長度)<sup>2</sup> ≤ C  
人為設定

- $\mathcal{H}(C)$ : overlaps but not exactly the same as  $\mathcal{H}'_2$
- soft and smooth structure over  $C \geq 0$ :

$$\mathcal{H}(0) \subset \mathcal{H}(1.126) \subset \dots \subset \mathcal{H}(1126) \subset \dots \subset \mathcal{H}(\infty) = \mathcal{H}_{10}$$

沒有 constraint.

regularized hypothesis  $\mathbf{w}_{\text{REG}}$   
optimal solution from  
regularized hypothesis set  $\mathcal{H}(C)$

Regularized.

**Questions?**

# Matrix Form of Regularized Regression Problem

$$\min_{\mathbf{w} \in \mathbb{R}^{Q+1}} E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{z}_n - y_n)^2$$

$(\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y}) \leftarrow \text{matrix form}$

$$\text{s.t.} \quad \sum_{q=0}^Q w_q^2 \leq C \leftarrow \text{加上 constraint 的 regression}$$

$\mathbf{w}^T \mathbf{w}$

- $\sum_n \dots = (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$ , remember? :-) 1个半径是 $\sqrt{C}$ 的球体
- $\mathbf{w}^T \mathbf{w} \leq C$ : feasible  $\mathbf{w}$  within a radius- $\sqrt{C}$  hypersphere

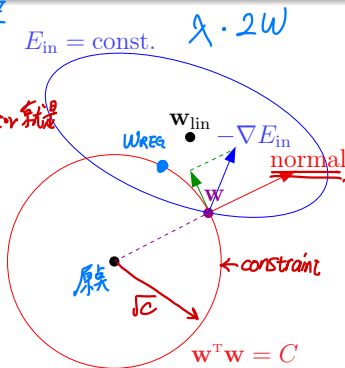
how to solve  
constrained optimization problem?

# The Lagrange Multiplier

$$\min_{\mathbf{w} \in \mathbb{R}^{Q+1}} E_{in}(\mathbf{w}) = \frac{1}{N}(\mathbf{Z}\mathbf{w} - \mathbf{y})^T(\mathbf{Z}\mathbf{w} - \mathbf{y}) \text{ s.t. } \mathbf{w}^T\mathbf{w} \leq C$$

$$\frac{\partial}{\partial \mathbf{w}} \downarrow \frac{2}{N}(\mathbf{Z}\mathbf{w} - \mathbf{y}) \cdot \mathbf{Z}$$

- decreasing direction:  $-\nabla E_{in}(\mathbf{w})$ ,  
**remember? :-)** 負 gradient
- normal vector of  $\mathbf{w}^T\mathbf{w} = C$ :  $\mathbf{w}$ . normal vector 就是  $\mathbf{w}$  本身
- if  $-\nabla E_{in}(\mathbf{w})$  and  $\mathbf{w}$  not parallel: can **decrease  $E_{in}(\mathbf{w})$  without violating the constraint** \*找  $-\nabla E_{in}(\mathbf{w})$  中垂直於  $\mathbf{w}$  的分量
- at optimal solution  $\mathbf{w}_{REG}$ ,  
 $-\nabla E_{in}(\mathbf{w}_{REG}) \propto \mathbf{w}_{REG}$  ← 最後 gradient 与  $\mathbf{w}$  平行



want: find Lagrange multiplier  $\lambda > 0$  and  $\mathbf{w}_{REG}$   
such that  $\nabla E_{in}(\mathbf{w}_{REG}) + \frac{2\lambda}{N} \mathbf{w}_{REG} = \mathbf{0}$

# Augmented Error

- if **oracle** tells you  $\lambda > 0$ , then

solving  $\nabla E_{\text{in}}(\mathbf{w}_{\text{REG}}) + \frac{2\lambda}{N} \mathbf{w}_{\text{REG}} = \mathbf{0}$

*Handwritten notes: 微分 (under the gradient), 2λ (circled), End. (above the gradient), ? (above the weight vector)*

$$\frac{2}{N} (Z^T Z \mathbf{w}_{\text{REG}} - Z^T \mathbf{y}) + \frac{2\lambda}{N} \mathbf{w}_{\text{REG}} = \mathbf{0}$$

- optimal solution:

*化簡 (simplified)*

*λ > 0 則 inverse 存在 (λ > 0 then inverse exists)*

$$\mathbf{w}_{\text{REG}} \leftarrow (Z^T Z + \lambda \mathbf{I})^{-1} Z^T \mathbf{y}$$

*(1 + λ)⁻¹ y = w*

—called ridge regression in Statistics

*linear regression 的進階版 (advanced version of linear regression)*

minimizing **unconstrained**  $E_{\text{aug}}$  effectively  
minimizes some **C-constrained**  $E_{\text{in}}$

# Augmented Error

- if **oracle** tells you  $\lambda > 0$ , then

solving

$$\nabla E_{in}(\mathbf{w}_{REG}) + \frac{2\lambda}{N} \boxed{\mathbf{w}_{REG}} = \mathbf{0}$$

(勸)

equivalent to minimizing

$$\underbrace{E_{in}(\mathbf{w})}_{\text{augmented error } E_{aug}(\mathbf{w})} + \underbrace{\frac{\lambda}{N} \mathbf{w}^T \mathbf{w}}_{\text{regularizer}} \leftarrow \text{最小化}$$

↑ 多一項

- regularization with (augmented error) instead of **constrained**  $E_{in}$

直接解  $E_{aug}(\mathbf{w})$

$$\mathbf{w}_{REG} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} E_{aug}(\mathbf{w}) \text{ for given } \lambda > 0 \text{ or } [\lambda = 0] \Rightarrow \text{沒有 constraint}$$

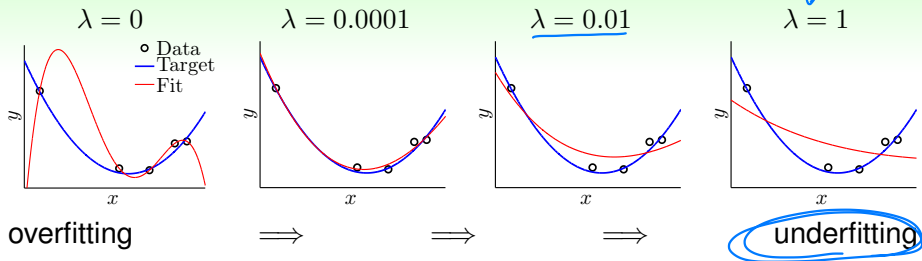
\*但要先給定  $\lambda$

minimizing unconstrained  $E_{aug}$  effectively  
minimizes some **C-constrained**  $E_{in}$



# The Results

constraint 太强  
↓



philosophy: a little regularization goes a long way!

call ' $+\frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$ ' weight-decay regularization:

長度

larger  $\lambda \leftarrow$  惩罚 long  $\mathbf{w}$

$\Leftrightarrow$  prefer shorter  $\mathbf{w}$

$\Leftrightarrow$  effectively smaller  $C$

—go with 'any transform + linear model'

**Questions?**

# Regularization and VC Theory

Regularization by  
Constrained-Minimizing  $E_{in}$

$$\min_{\mathbf{w}} E_{in}(\mathbf{w}) \text{ s.t. } \mathbf{w}^T \mathbf{w} \leq C$$



VC Guarantee of  
Constrained-Minimizing  $E_{in}$

$$E_{out}(\mathbf{w}) \leq E_{in}(\mathbf{w}) + \Omega(\underline{\underline{\mathcal{H}(C)}})$$



C equivalent to some  $\lambda$  給定C等於給入

Regularization by  
Minimizing  $E_{aug}$

$$\min_{\mathbf{w}} E_{aug}(\mathbf{w}) = E_{in}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

因為是 unconstrained, 所以還是考慮了所有  $\mathbf{w}$   
只是不去選它們

minimizing  $E_{aug}$ : indirectly getting VC  
guarantee without confining to  $\mathcal{H}(C)$

# Another View of Augmented Error

## Augmented Error

$$E_{\text{aug}}(\mathbf{w}) \stackrel{=}{=} E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

## VC Bound

$$E_{\text{out}}(\mathbf{w}) \stackrel{\leq}{=} E_{\text{in}}(\mathbf{w}) + \underbrace{\Omega(\mathcal{H})}_{\text{complexity penalty}}$$

- regularizer  $\mathbf{w}^T \mathbf{w}$  : complexity of a single hypothesis
- generalization price  $\Omega(\mathcal{H})$ : complexity of a hypothesis set
- if  $\frac{\lambda}{N} \Omega(\mathbf{w})$  'represents'  $\Omega(\mathcal{H})$  well,  $E_{\text{aug}}$  is a better proxy of  $E_{\text{out}}$  than  $E_{\text{in}}$

$\frac{\lambda}{N} \Omega(\mathbf{w})$  的 complexity

代理人

minimizing  $E_{\text{aug}}$ :

(heuristically) operating with the better proxy;  
(technically) enjoying flexibility of whole  $\mathcal{H}$

# Effective VC Dimension

$$\min_{\mathbf{w} \in \mathbb{R}^{\tilde{d}+1}} E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \Omega(\mathbf{w})$$

- model complexity? ← 沒變  
 $d_{\text{VC}}(\mathcal{H}) = \tilde{d} + 1$ , because  $\{\mathbf{w}\}$  'all considered' during minimization
- $\{\mathbf{w}\}$  'actually needed':  $\mathcal{H}(\mathcal{C})$ , with some  $\mathcal{C}$  equivalent to  $\lambda$   
 ↑  $\mathcal{H}(\mathcal{C})$  內的  $\mathbf{w}$  才是真的會被考慮的
- $d_{\text{VC}}(\mathcal{H}(\mathcal{C}))$ :  
effective VC dimension  $d_{\text{EFF}}(\mathcal{H}, \underbrace{\mathcal{A}}_{\text{algorithm}})$   
 等效 VC dimension  $\min E_{\text{aug}}$

explanation of regularization:

$d_{\text{VC}}(\mathcal{H})$  large,  
 while  $d_{\text{EFF}}(\mathcal{H}, \mathcal{A})$  small if  $\mathcal{A}$  regularized

# Questions?

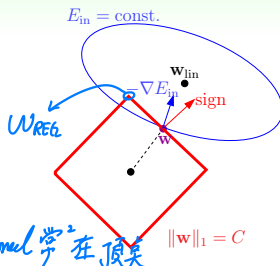
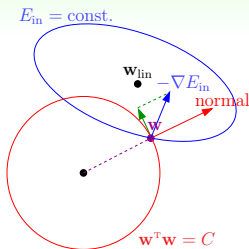
# General Regularizers $\Omega(\mathbf{w})$

want: constraint in the 'direction' of target function

- target-dependent: some **properties** of target, if known
  - symmetry regularizer:  $\sum \mathbb{I}[q \text{ is odd}] w_q^2 \leftarrow$  对称: 让奇数次方项变小
- plausible: direction towards **smoother** or **simpler**.  
 stochastic/deterministic noise both **non-smooth**
  - sparsity (L1) regularizer:  $\sum |w_q|$  (next slide)  $\leftarrow$  找出简单的
- friendly: easy to **optimize**
  - weight-decay** (L2) regularizer:  $\sum w_q^2 \leftarrow$  容易 Optimize.
- bad? :-)**: no worries, guard by  $\lambda$ .

augmented error = error  $\widehat{\text{err}}$  + regularizer  $\Omega$   
 regularizer: **target-dependent**, **plausible**, or **friendly**  
*domain knowledge*

# L2 and L1 Regularizer



## L2 Regularizer

$$\Omega(\mathbf{w}) = \sum_{q=0}^Q w_q^2 = \|\mathbf{w}\|_2^2$$

- convex, differentiable everywhere
- easy to optimize

## L1 Regularizer

$$\Omega(\mathbf{w}) = \sum_{q=0}^Q |w_q| = \|\mathbf{w}\|_1$$

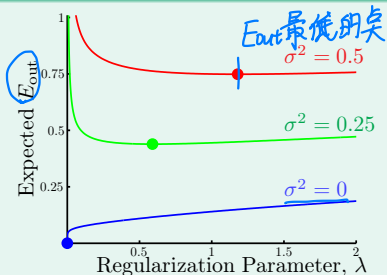
- convex, not differentiable everywhere
  - sparsity in solution
- 大部分 weight 都是零

L1 useful if needing sparse solution. 省計算

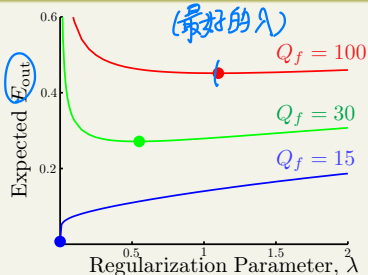


# The Optimal $\lambda$

## stochastic noise



## deterministic noise



- more noise  $\iff$  more regularization needed  
—more bumpy road  $\iff$  putting brakes more
- noise **unknown**—important to **make proper choices**

how to choose?

stay tuned for the next lecture! :-)

**Questions?**

# Summary

## 1 How Can Machines Learn?

### Lecture 06: Beyond Basic Linear Models

## 2 How Can Machines Learn **Better**?

### Lecture 07: Combatting Overfitting

- What is Overfitting?  
**lower  $E_{in}$  but higher  $E_{out}$**
- The Role of Noise and Data Size  
**overfitting 'easily' happens!**
- Deterministic Noise  
**what  $\mathcal{H}$  cannot capture acts like noise**
- Dealing with Overfitting  
**data cleaning/pruning/hinting & regularization**
- Regularized Hypothesis Set  
**original  $\mathcal{H}$  + constraint**
- Weight Decay Regularization  
**add  $\frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$  in  $E_{aug}$**
- Regularization and VC Theory  
**regularization decreases  $d_{EFF}$**
- General Regularizers  
**target-dependent, [plausible], or [friendly]**

- **next: choosing from the so-many models/parameters**