

Machine Learning

(機器學習)

Lecture 8: Combatting Overfitting (2)

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 8: Combatting Overfitting (2)

- Model Selection Problem
- Validation
- Leave-One-Out Cross Validation
- V-Fold Cross Validation

So Many Models Learned

Even Just for Binary Classification . . .

$\mathcal{A} \in \{ \text{PLA, pocket, linear regression, logistic regression} \}$

×

$T \in \{ 100, 1000, 10000 \} \leftarrow \text{走幾步}$

×

$\eta \in \{ 1, 0.01, 0.0001 \} \leftarrow \text{learning rate}$

×

$\Phi \in \{ \text{linear, quadratic, poly-10, Legendre-poly-10} \} \leftarrow \text{feature transform}$

×

$\Omega(\mathbf{w}) \in \{ \text{L2 regularizer, L1 regularizer, symmetry regularizer} \}$

×

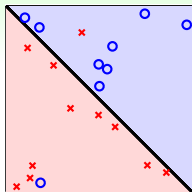
$\lambda \in \{ 0, 0.01, 1 \} \leftarrow \lambda$

↑
regularizer?

太多組合

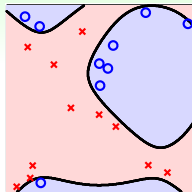
in addition to your **favorite** combination, may need to try other combinations to get a good g

Model Selection Problem



\mathcal{H}_1

which one do you prefer? :-)



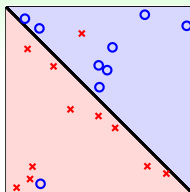
\mathcal{H}_2

- given: M models $\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_M$, each with corresponding algorithm $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_M$. *选出一个最佳 model: \mathcal{H}_{m^*}*
- goal: select \mathcal{H}_{m^*} such that $g_{m^*} = \mathcal{A}_{m^*}(\mathcal{D})$ is of low $E_{\text{out}}(g_{m^*})$
- unknown E_{out} due to unknown $P(\mathbf{x})$ & $P(y|\mathbf{x})$, as always :-)
- arguably the most important practical problem of ML

↓ *如何选 model*

how to select? **visually?**
—no, remember Lecture 7? :-)

Model Selection by Best E_{in}

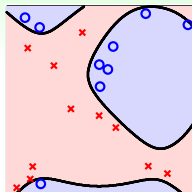


\mathcal{H}_1

选更複雜

select by best E_{in} ?

$$\underline{m^*} = \underset{1 \leq m \leq M}{\operatorname{argmin}} (E_m = E_{in}(\mathcal{A}_m(\mathcal{D})))$$



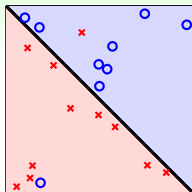
\mathcal{H}_2

- Φ_{1126} always more preferred over Φ_1 ; $\lambda = 0$ always more preferred over $\lambda = 0.1$ — **overfitting?**
- if \mathcal{A}_1 minimizes E_{in} over \mathcal{H}_1 and \mathcal{A}_2 minimizes E_{in} over \mathcal{H}_2 ,
 $\implies g_{m^*}$ achieves minimal E_{in} over $\mathcal{H}_1 \cup \mathcal{H}_2$ \implies 等於在考慮 $\mathcal{H}_1 \cup \mathcal{H}_2$ 的 hypothesis
 \implies 'model selection + learning' pays $d_{vc}(\mathcal{H}_1 \cup \mathcal{H}_2)$ hypothesis
 — **bad generalization?**

用 E_{in} 选很危險

selecting by E_{in} is **dangerous**

Model Selection by Best E_{test}

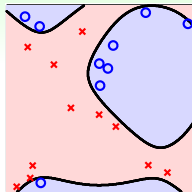


\mathcal{H}_1

select by best E_{test} , which is evaluated on a fresh $\mathcal{D}_{\text{test}}$?

$$m^* = \operatorname{argmin}_{1 \leq m \leq M} (E_m = E_{\text{test}}(\mathcal{A}_m(\mathcal{D})))$$

选 E_{test} 最小的 model



\mathcal{H}_2

- generalization guarantee (finite-bin Hoeffding):

$$E_{\text{out}}(g_{m^*}) \leq E_{\text{test}}(g_{m^*}) + O\left(\sqrt{\frac{\log M}{N_{\text{test}}}}\right)$$

—yes! strong guarantee :-)

E_{out} 被 bound 住。

- but where is $\mathcal{D}_{\text{test}}$? —your boss's safe, maybe? :-)

我得到 $\mathcal{D}_{\text{test}}$?

selecting by E_{test} is **infeasible** and **cheating**

Comparison between E_{in} and E_{test}

in-sample error E_{in}

- calculated from \mathcal{D}
- **feasible** on hand
- 'contaminated' as \mathcal{D} also used by \mathcal{A}_m to 'select' g_m

test error E_{test}

- calculated from \mathcal{D}_{test}
- infeasible in boss's safe
- 'clean' as \mathcal{D}_{test} never used for selection before

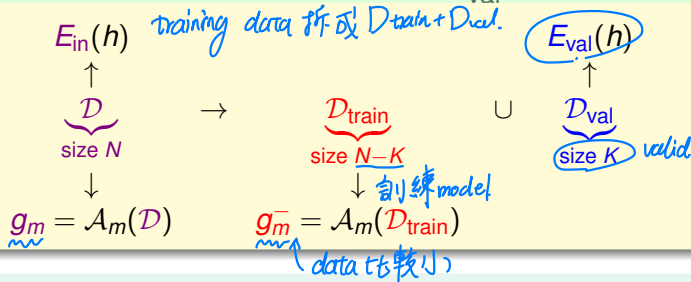
something in between: E_{val}

- calculated from $\mathcal{D}_{val} \subset \mathcal{D}$
- **feasible** on hand *validation set*
- 'clean' **if** \mathcal{D}_{val} never used by \mathcal{A}_m before

selecting by E_{val} : legal cheating :-)

Questions?

Validation Set \mathcal{D}_{val}



- $\mathcal{D}_{\text{val}} \subset \mathcal{D}$: called **validation set**—‘on-hand’ simulation of test set
- to connect E_{val} with E_{out} : $\mathcal{D}_{\text{val}} \stackrel{\text{iid}}{\sim} P(\mathbf{x}, y) \iff$ select K examples from \mathcal{D} at random (blue arrow: \mathcal{D}_{val} 從 \mathcal{D} 平均隨機抽樣)
- to make sure \mathcal{D}_{val} ‘clean’:
feed only $\mathcal{D}_{\text{train}}$ to \mathcal{A}_m for model selection

$$E_{\text{out}}(\underline{g_m}) \leq \underline{E_{\text{val}}}(\underline{g_m}) + O\left(\sqrt{\frac{\log M}{K}}\right)$$

Model Selection by Best E_{val}

選 val error 最小的 model

$$m^* = \underset{1 \leq m \leq M}{\operatorname{argmin}} (E_m = E_{\text{val}}(\mathcal{A}_m(\mathcal{D}_{\text{train}})))$$

- generalization guarantee for all m :

$$E_{\text{out}}(\underline{g_m}) \leq E_{\text{val}}(\underline{g_m}) + O\left(\sqrt{\frac{\log M}{K}}\right)$$

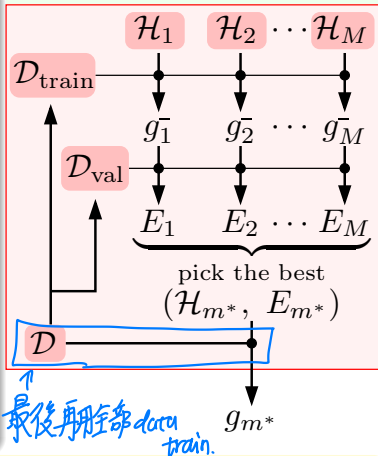
- heuristic gain from $N - K$ to N :

選完 model 之後, 連同 validation set 一起下去 train.

$$E_{\text{out}}\left(\underbrace{g_{m^*}}_{\mathcal{A}_{m^*}(\mathcal{D})}\right) \leq E_{\text{out}}\left(\underbrace{g_{m^*}^-}_{\mathcal{A}_{m^*}(\mathcal{D}_{\text{train}})}\right)$$

↑ data 又多了一點

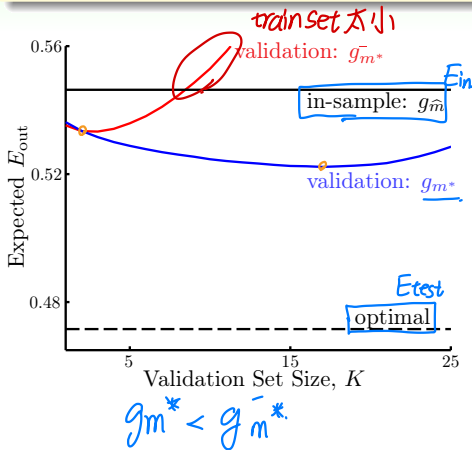
— learning curve, remember? :-)



$$\underline{E_{\text{out}}}(g_{m^*}) \leq E_{\text{out}}(g_{m^*}^-) \leq E_{\text{val}}(g_{m^*}^-) + O\left(\sqrt{\frac{\log M}{K}}\right)$$

Validation in Practice

use validation to select between $\mathcal{H}_{\Phi_5}^{5 \times 2}$ and $\mathcal{H}_{\Phi_{10}}^{10 \times 2}$



- in-sample: selection with E_{in}
- optimal: cheating-selection with E_{test}
- sub-g: selection with E_{val} and report \bar{g}_m^*
- full-g: selection with E_{val} and report \underline{g}_m^*
 — $E_{\text{out}}(\underline{g}_m^*) \leq E_{\text{out}}(\bar{g}_m^*)$
 indeed

why is sub-g worse than in-sample some time?

The Dilemma about K validation set 大小

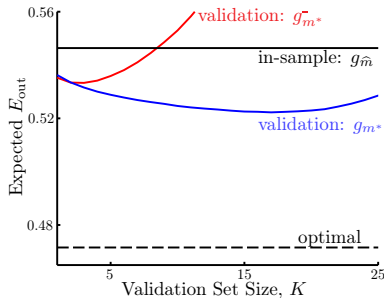
reasoning of validation:

大量 valid set 才能保證 E_{out} 与 E_{val} 接近

$$E_{out}(g) \underset{\text{(small } K)}{\approx} E_{out}(g^-) \underset{\text{(large } K)}{\approx} E_{val}(g^-)$$

但 valid set 拿太多, 又会說 training 時偏差方大

- large K : **every** $E_{val} \approx E_{out}$,
but all g_m^- much worse than g_m
- small K : every $g_m^- \approx g_m$,
but E_{val} far from E_{out}



practical rule of thumb: $K = \frac{N}{5}$

20%

Questions?

Extreme Case: $K = 1$ K 超小

reasoning of validation: ✓

$$E_{\text{out}}(g) \approx E_{\text{out}}(g^-) \approx E_{\text{val}}(g^-)$$

(small K) (large K)

GG!

- take $K = 1$? $\mathcal{D}_{\text{val}}^{(n)} = \{(\mathbf{x}_n, y_n)\}$ and $E_{\text{val}}^{(n)}(g_n^-) = \text{err}(g_n^-(\mathbf{x}_n), y_n) = e_n$ 第 N 筆資料的 error
- make e_n closer to $E_{\text{out}}(g)$?—average over possible $E_{\text{val}}^{(n)}$
- leave-one-out cross validation estimate:

$$E_{\text{loocv}}(\mathcal{H}, \mathcal{A}) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N \text{err}(g_n^-(\mathbf{x}_n), y_n)$$

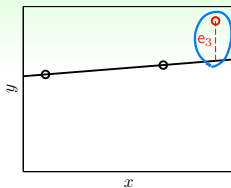
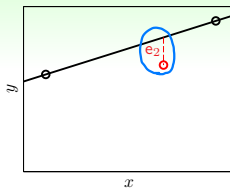
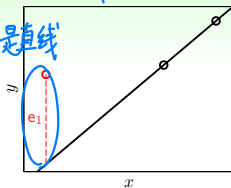
leave-one-out cross validation. 平均

hope: $E_{\text{loocv}}(\mathcal{H}, \mathcal{A}) \approx E_{\text{out}}(g)$ 希望

Illustration of Leave-One-Out

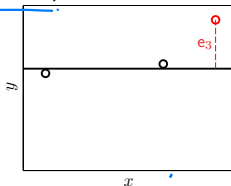
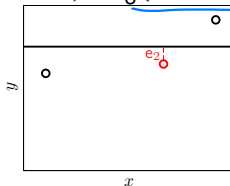
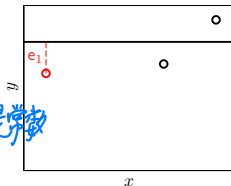
假設只有3行 data point

hypothesis是直線



$$E_{\text{loocv}}(\text{linear}) = \frac{1}{3}(e_1 + e_2 + e_3)$$

hypothesis是常數



$$E_{\text{loocv}}(\text{constant}) = \frac{1}{3}(e_1 + e_2 + e_3)$$

用常數 model 就夠了

which one would you choose?

$$m^* = \underset{1 \leq m \leq M}{\operatorname{argmin}} (E_m = E_{\text{loocv}}(\mathcal{H}_m, \mathcal{A}_m))$$

Theoretical Guarantee of Leave-One-Out Estimate

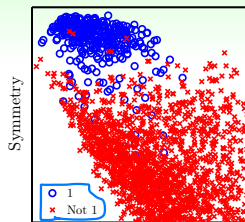
does $E_{\text{loocv}}(\mathcal{H}, \mathcal{A})$ say something about $E_{\text{out}}(g)$?

yes, for average E_{out} on size- $(N - 1)$ data

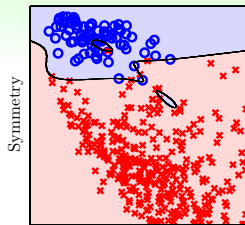
$$\begin{aligned}
 \underbrace{\mathcal{E}_{\mathcal{D}}}_{\text{期望值}} E_{\text{loocv}}(\mathcal{H}, \mathcal{A}) &= \underbrace{\mathcal{E}_{\mathcal{D}} \frac{1}{N} \sum_{n=1}^N e_n}_{\text{期望值}} = \frac{1}{N} \sum_{n=1}^N \mathcal{E}_{\mathcal{D}} e_n \\
 &= \frac{1}{N} \sum_{n=1}^N \mathcal{E}_{\mathcal{D}_n(\mathbf{x}_n, y_n)} \left\{ \mathcal{E} \left[\text{err}(g_n^-(\mathbf{x}_n), y_n) \right] \right\} \quad \text{leave one-out 的那串 data} \\
 &= \frac{1}{N} \sum_{n=1}^N \mathcal{E}_{\mathcal{D}_n} \left[E_{\text{out}}(g_n^-) \right] \\
 &= \frac{1}{N} \sum_{n=1}^N \overline{E_{\text{out}}(N-1)} = \underbrace{\overline{E_{\text{out}}(N-1)}}_{\text{平均值}} = \underbrace{\overline{E_{\text{out}}(g^-)}}_{E_{\text{out}}(g)}
 \end{aligned}$$

expected $E_{\text{loocv}}(\mathcal{H}, \mathcal{A})$ says something about expected $E_{\text{out}}(g^-)$
 —often called 'almost unbiased estimate of $E_{\text{out}}(g)$ '

Leave-One-Out in Practice

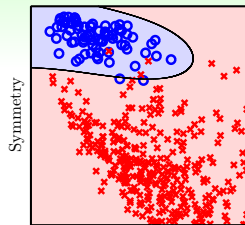


1 or not 1



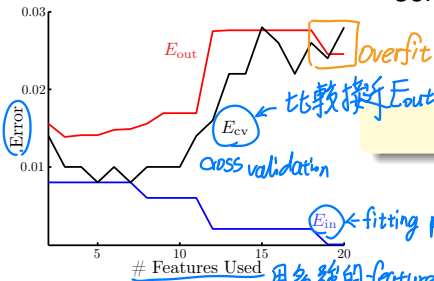
Average Intensity

select by E_{in}



Average Intensity

select by E_{loocv}



E_{loocv} much better than E_{in}

用多強的 feature transformation.

Questions?

Disadvantages of Leave-One-Out Estimate

Computation

$$E_{\text{loocv}}(\mathcal{H}, \mathcal{A}) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N \text{err}(\mathbf{g}_n^-(\mathbf{x}_n), y_n)$$

- N 'additional' training per model, not always feasible in practice
- except 'special case' like analytic solution for linear regression

linear regression 適合用 LOOCV

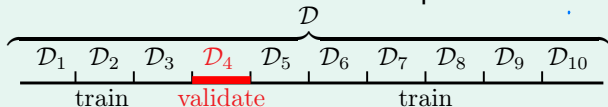
E_{loocv} : not often used practically

不好搞

V-fold Cross Validation

how to **decrease computation need** for cross validation?

- essence of leave-one-out cross validation: partition \mathcal{D} to N parts, taking $N - 1$ for training and 1 for validation orderly
- V-fold cross-validation: random-partition of \mathcal{D} to **V equal parts**,



只拿一份V, 做 validation

take $V - 1$ for training and 1 for validation orderly

$$E_{cv}(\mathcal{H}, \mathcal{A}) = \frac{1}{V} \sum_{v=1}^V \underline{E_{val}^{(v)}(g_v^-)}$$

平均

- selection by E_{cv} : $\underline{m^*} = \underset{1 \leq m \leq M}{\operatorname{argmin}} (E_m = \underline{E_{cv}}(\mathcal{H}_m, \mathcal{A}_m))$
cross validation.

practical rule of thumb:

V = 10

20%? 10%?

Final Words on Validation

'Selecting' Validation Tool

- **V-Fold** generally preferred over single validation if computation allows
cross validation.
- **5-Fold or 10-Fold** generally works well: *通常不用 200*
not necessary to trade V-Fold with Leave-One-Out

Nature of Validation

- all training models: select among hypotheses
- all validation schemes: **select among finalists**
- all testing methods: just evaluate

validation still **more optimistic than testing**

do not fool yourself and others :-),
report test result, not **best validation result**

Questions?

Summary

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

Lecture 8: Combatting Overfitting (2)

- Model Selection Problem
dangerous by E_{in} and dishonest by E_{test}
- Validation
select with $E_{\text{val}}(\mathcal{A}_m(\mathcal{D}_{\text{train}}))$ while returning $\mathcal{A}_{m^*}(\mathcal{D})$
- Leave-One-Out Cross Validation
huge computation for almost unbiased estimate
- V-Fold Cross Validation
reasonable computation and performance
- **next: something 'up my sleeve'**