

4/2 月日



M2608.001300

Machine Learning Fundamentals & Applications

[2: Feasibility of Learning]

Electrical and Computer Engineering
Seoul National University

© 2018 Sungroh Yoon. this material is for educational uses only. some contents are based on
the material provided by the textbook authors and may be copyrighted by them.

Outline

Learning Unknown Target Function

- Hoeffding Inequality

- Connection to Learning

- Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

- Error Measures

- Noisy Targets

Summary

Readings

- *Learning from Data* by Abu-Mostafa, Magdon-Ismail, and Lin
 - ▷ Chapter 1: The Learning Problem (Sections 1.3 & 1.4)

Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

Error Measures

Noisy Targets

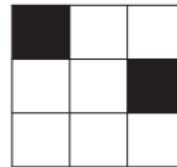
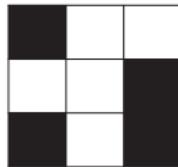
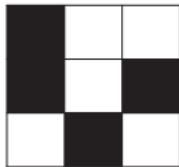
Summary

Learning the target function

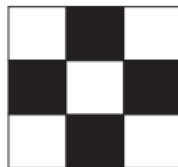
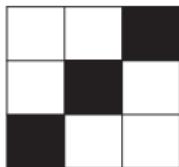
- target function f = learning objective
 - ▶ most important fact about f : it is unknown
- natural question
 - ▶ how could a limited data set reveal enough information to determine the entire f ?



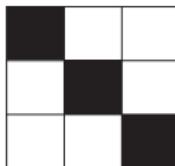
- 6 training examples of a binary target function f



$$f = -1$$



$$f = +1$$



$$f = ?$$

- difficulty in this simple problem: the rule, not the exception

Inference outside training data

- in training data \mathcal{D} , we know the value of f on all points
 - ▶ but this does not mean that we have learned f
- key question: does \mathcal{D} tell us anything outside \mathcal{D} we did not know before?
 - ▶ if so, then we have learned *something*
 - ▶ otherwise, we can conclude: learning is not feasible

Example

- consider a Boolean f over 3D input space $\mathcal{X} = \{0, 1\}^3$
 - ▶ # distinct input vectors: $2^3 = 8$
 - ▶ # distinct Boolean functions on 3 Boolean inputs: $2^{2^3} = 2^8 = 256$

- given: a data set \mathcal{D} of five examples:

x_n	y_n
000	○
001	●
010	●
011	○
100	●

- task: predict f on points never seen before
= outside \mathcal{D}

- any function that agrees with \mathcal{D} could be f (e.g. f_1, \dots, f_8 below)

x	y	g	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8
000	○	○	○	○	○	○	○	○	○	○
001	●	●	●	●	●	●	●	●	●	●
010	●	●	●	●	●	●	●	●	●	●
011	○	○	○	○	○	○	○	○	○	○
100	●	●	●	●	●	●	●	●	●	●
101	?	○	○	○	○	●	●	●	●	●
110	?	○	○	●	●	○	○	●	●	●
111	?	○	●	○	●	○	●	○	●	●

- purpose of learning

▶ find such g that agrees the most with f

- candidates of g

- ▶ g_1 : returns \bullet for all the three points
- ▶ g_2 : returns \circ for all the three points
- ▶ g_3 : returns $\text{XOR}(x)^1$ for all the three points
- ▶ g_4 : returns the opposite of g_3 for all the three points

x	y	g	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	g_1	g_2	g_3	g_4
000	○	○	○	○	○	○	○	○	○	○				
001	●	●	●	●	●	●	●	●	●	●	●	●	●	●
010	●	●	●	●	●	●	●	●	●	●	●	●	●	●
011	○	○	○	○	○	○	○	○	○	○	○	○	○	○
100	●	●	●	●	●	●	●	●	●	●	●	●	●	●
101		?	○	○	○	○	●	●	●	●	●	●	●	●
110		?	○	○	●	●	○	○	●	●	●	●	●	●
111		?	○	●	○	●	○	●	○	●	●	●	●	●

- which is the best?

¹ $\text{XOR}(x)$ returns \bullet if # 1's in x is odd; returns \circ otherwise

x	y	g	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	g_1	g_2	g_3	g_4
000	○	○	○	○	○	○	○	○	○	○				
001	●	●	●	●	●	●	●	●	●	●				
010	●	●	●	●	●	●	●	●	●	●				
011	○	○	○	○	○	○	○	○	○	○				
100	●	●	●	●	●	●	●	●	●	●				
101		?	○	○	○	○	●	●	●	●	●	○	○	●
110		?	○	○	●	●	○	○	●	●	●	○	○	●
111		?	○	●	○	●	○	●	○	●	●	○	●	○

- performance evaluation

- ▶ in terms of # points in agreement between g_k and f_i

g_k	on three	on two	on one	on none
g_1	f_8	f_4, f_6, f_7	f_2, f_3, f_5	f_1
g_2	f_1	f_2, f_3, f_5	f_4, f_6, f_7	f_8
g_3	f_2	$\cancel{f_1} \cancel{f_4} \cancel{f_6}$	$\cancel{f_3} \cancel{f_5} \cancel{f_8}$	f_7
g_4	f_7	f_3, f_5, f_8	f_1, f_4, f_6	f_2

Are we doomed?

- no matter what \mathcal{A} does or what \mathcal{H} is used
 - ▶ no difference as far as performance outside \mathcal{D} is concerned
all that matters in learning
- learning an unknown function f ? \Rightarrow impossible!
 - ▶ f can assume any value outside \mathcal{D}
 - ▶ we can even prove: f remains unknown outside \mathcal{D}
- so what now?

probability를 이용해
 f 와 학습집단을 유사한
 g 를 학습하는 것이 가능하다.



Outline

Learning Unknown Target Function

- Hoeffding Inequality

- Connection to Learning

- Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

- Error Measures

- Noisy Targets

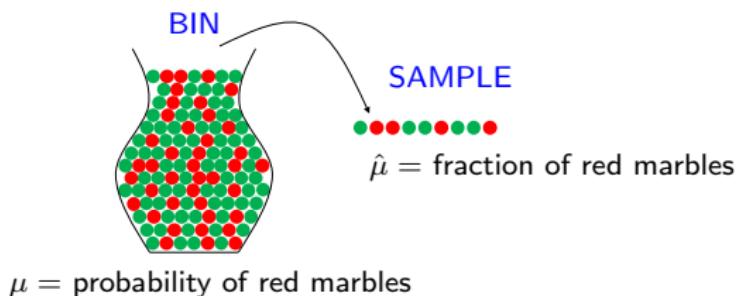
Summary

Rescue plan

- probability
- we can infer something outside \mathcal{D} using only \mathcal{D}
 - ▶ in a probabilistic way
 - ▶ we give up knowing about f “for sure”

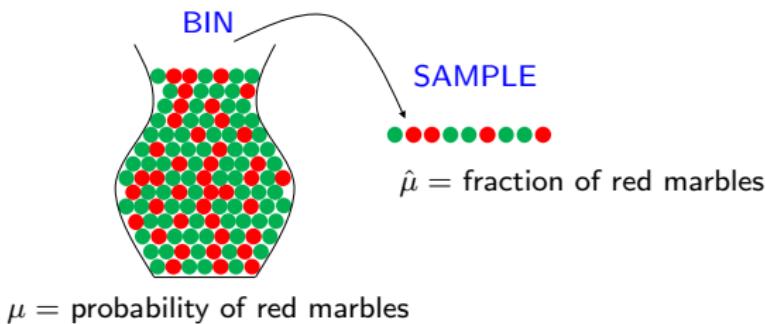
A 'bin' experiment

- consider a bin with red and green marbles (possibly infinitely many)
 $\mathbb{P}[\text{picking a red marble}] = \mu, \quad \mathbb{P}[\text{picking a green marble}] = 1 - \mu$
 - ▶ the value of μ : Unknown
- pick N marbles ~~independently~~ (with replacement)
 - ▶ the fraction of red marbles within the sample $\triangleq \hat{\mu}$



Does $\hat{\mu}$ say anything about μ ?

- **no** e.g. sample can be mostly green while bin is mostly red
 - ▶ *possible but not probable*
- **yes** sample frequency $\hat{\mu}$ is likely 逼近 bin frequency μ



- how to quantify the relationship between $\hat{\mu}$ and μ ?

Hoeffding inequality

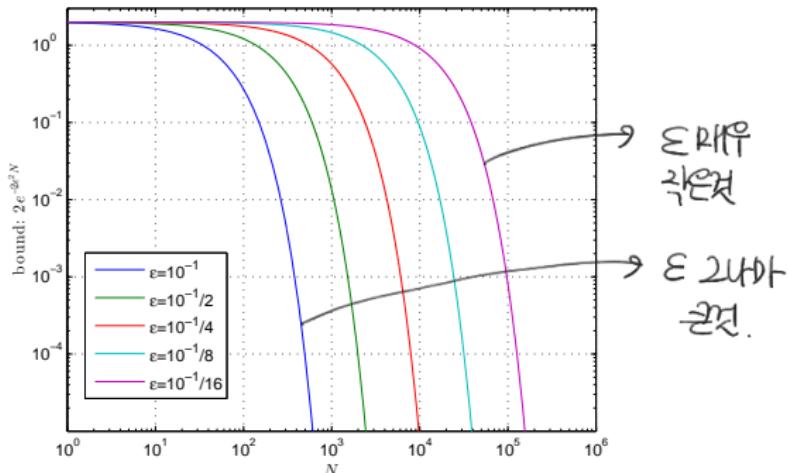
for any sample size N and $\epsilon > 0$:

$$\boxed{\mathbb{P}[|\hat{\mu} - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}} \quad (1)$$

- ▶ \mathbb{P} : probability of an event (here, wrt random sample)
- ▶ the bound does not depend on μ or on bin size
- in words
 - ▶ as the sample size N grows, it becomes exponentially unlikely that $\hat{\mu}$ will deviate from μ by more than our ‘tolerance’ ϵ
- in other words, the statement “ $\mu = \hat{\mu}$ ” is
 - ▶ PAC (probably approximately correct)

$$\mathbb{P}[|\hat{\mu} - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N} \quad (1)$$

- if we choose ϵ to be very small (to make $\hat{\mu}$ approximate μ well)
 - ▶ need a $b_i g$ sample size N to make the rhs of (1) small



Wassily Hoeffding

- the key player in the bound $2e^{-2\epsilon^2 N}$: N
 - ▶ if $N \rightarrow \infty$, $\mu \approx \hat{\mu}$ with very high probability (albeit *not for sure*)

Example

$$\mathbb{P}[|\hat{\mu} - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N} \quad (1)$$

- $N = 1,000$

99% of the time $\mu - 0.05 \leq \hat{\mu} \leq \mu + 0.05$ ($\epsilon = 0.05$)

99.9999996% of the time $\mu - 0.10 \leq \hat{\mu} \leq \mu + 0.10$ ($\epsilon = 0.10$)

- meaning: if we repeatedly pick a sample of size 1,000

$$\mu \in [\hat{\mu} - 0.05, \hat{\mu} + 0.05]$$

- ▶ and we are right 99% of the time about this
- we learned *something* : we reached outside \mathcal{D} to μ

Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

Error Measures

Noisy Targets

Summary

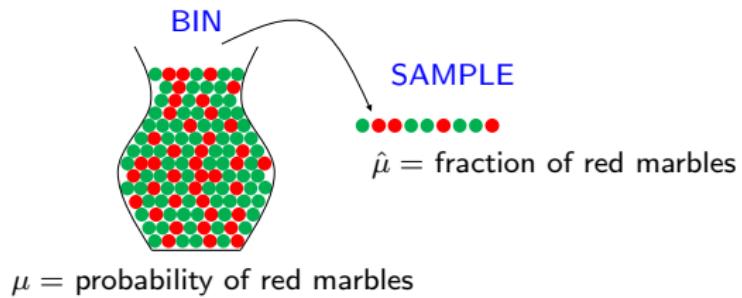
Are the bin model and learning connected?

- seemingly different
 - ▶ unknown in the bin problem: just the value of μ
 - ▶ unknown in learning: an entire function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- but we can connect them:
 - ▶ $\mu = \text{error rate}$ a hypothesis makes when approximating f



- obvious
 - ▶ marble $\equiv \mathbf{x} \in \mathcal{X}$
 - ▶ bin $\equiv \mathcal{X}$
 - ▶ sample $\equiv \mathcal{D}$

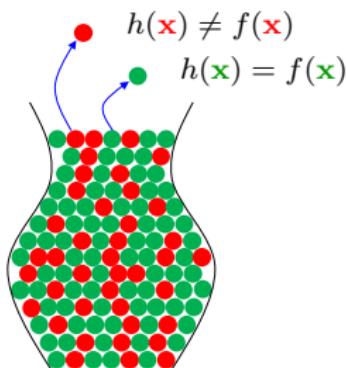
- less obvious
 - ▶ ball color $\equiv ?$
 - ▶ $\mu \equiv ?$
 - ▶ $\hat{\mu} \equiv ?$



Bin coloring \equiv hypothesis h

- take a *particular* hypothesis $h \in \mathcal{H}$ and compare it to f on each $\mathbf{x} \in \mathcal{X}$:

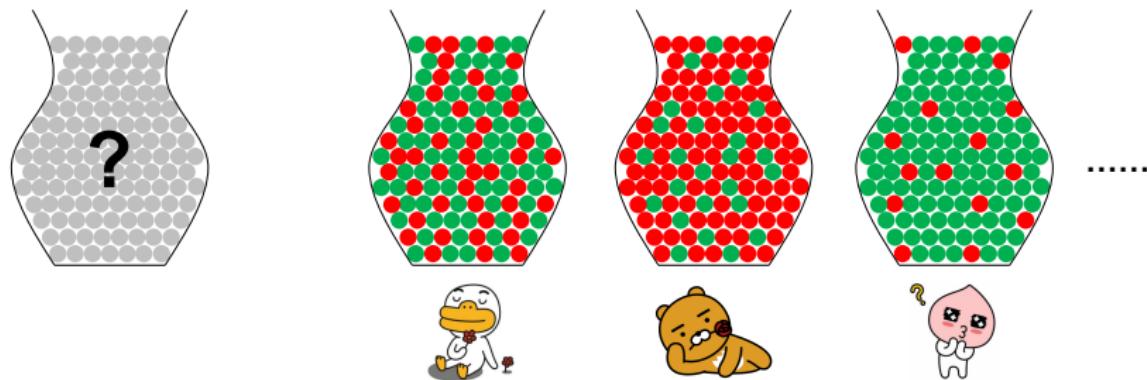
- ▶ color the point \mathbf{x} $\begin{cases} \text{green} & \text{if } h(\mathbf{x}) = f(\mathbf{x}) : \text{hypothesis right} \\ \text{red} & \text{if } h(\mathbf{x}) \neq f(\mathbf{x}) : \text{hypothesis wrong} \end{cases}$



- this means:

- ▶ $\mu \equiv$ error rate a hypothesis makes when approximating f
즉, 통계학적 연결성을 우리의 g 를 평가할 수 있다.

- different h
 - ⇒ different coloring of the balls in the bin
 - ⇒ different μ (*i.e.* error rate depends on the hypothesis used)



- learning algorithm \mathcal{A} picks the best h (*i.e.* g)

Sample \equiv training data set \mathcal{D}

- points $\mathbf{x}_1, \dots, \mathbf{x}_N$ in \mathcal{D}
 - ▶ picked independently according to distribution P
 - ▶ a random sample of red and green points
 - ▶ the color of each point in \mathcal{D} will be known to us
- P : some probability distribution over \mathcal{X} P will tell D to pick.
 - ▶ μ is allowed to be unknown $\rightarrow P$ can be unknown as well
- $\hat{\mu}$ estimates μ , telling us something about f
 - i.e. if $\hat{\mu} \approx 0 \Rightarrow h \approx f$ over entire \mathcal{X}

Revised: the basic learning step

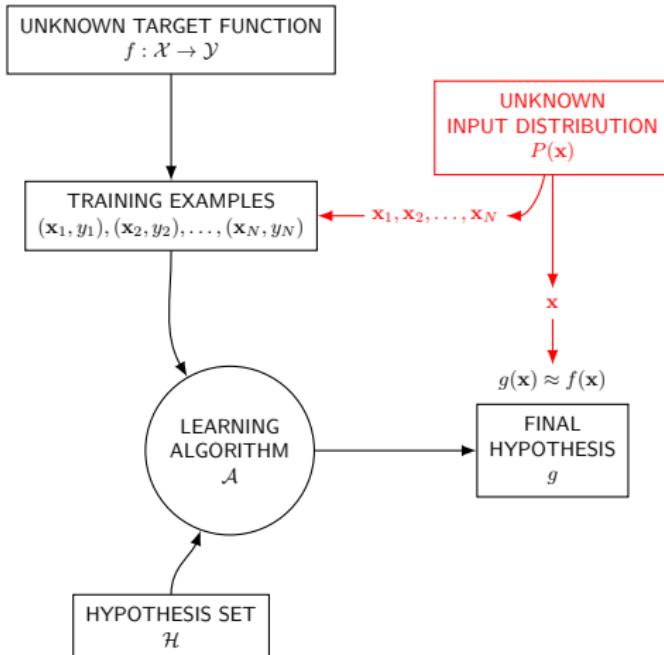


Figure 1: probabilistic component added to the basic learning step

Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound

$$P\{|E_m(g) - E_{out}(g)| > \varepsilon\} < \text{Hoeffding Bound}$$

Feasibility of Learning

할수있다.

Error and Noise

Error Measures

Noisy Targets

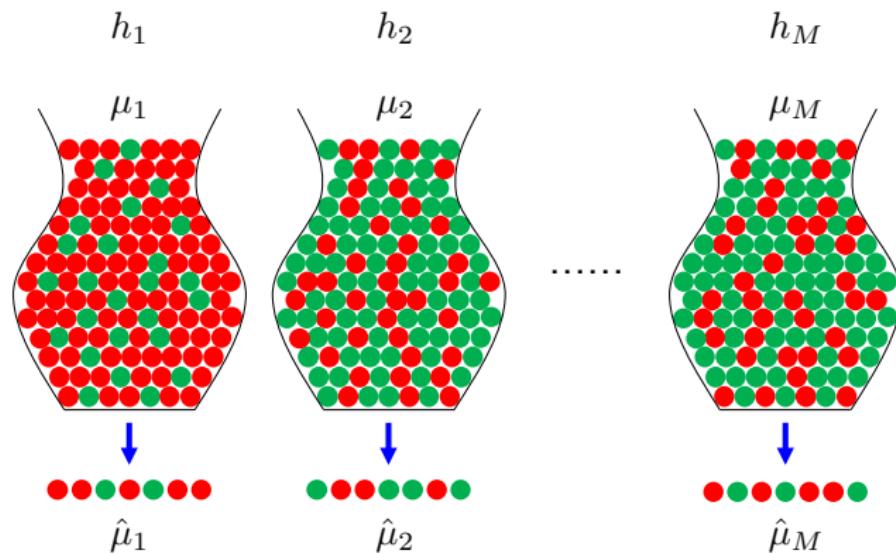
Summary

Are we done?

- no. how do we know $\hat{\mu} \approx 0$ in general?
 - in real learning
 - ▶ explore an entire hypothesis set \mathcal{H}
 - ▶ look for some $h \in \mathcal{H}$ that has a small error rate
 - ▶ choose the best from multiple h 's
- for example, "selection bias ~~term~~" is
one of the challenging problem.

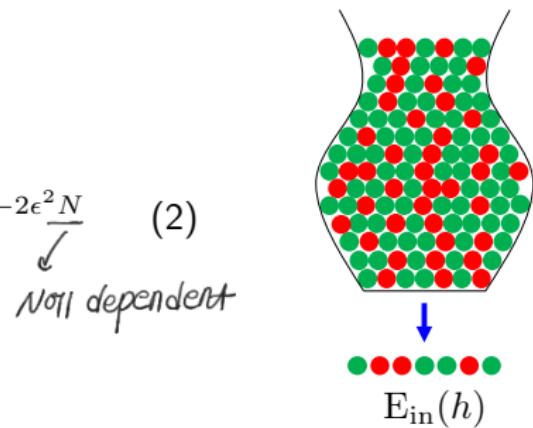
Multiple bins

- in order to capture real learning
 - ▶ extend the bin equivalence to multiple hypotheses
 - ▶ different $h \rightarrow$ different $\hat{\mu} \rightarrow$ different approximation of μ



- both $\hat{\mu}$ and our belief in μ depend on which hypothesis h is used
 - $\hat{\mu}$ is in sample error denoted by $E_{in}(h)$
 - μ is out of sample error denoted by $E_{out}(h)$

$$E_{out}(h)$$



- Hoeffding inequality now becomes

$$\mathbb{P}[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N} \quad (2)$$

Non dependent

In-sample error $E_{in}(h)$

- error rate within a sample:

$$E_{in}(h) = \text{(fraction of } D \text{ where } f \text{ and } h \text{ disagree)}$$
$$= \frac{1}{N} \sum_{n=1}^N [h(\mathbf{x}_n) \neq f(\mathbf{x}_n)]$$

(3)

\Rightarrow binary classification
의 예시
(multiple 은 다른 것
사용할 것)

▶ corresponds to $\hat{\mu}$ in the bin model

- ▶ notation: $[\text{true}] = 1, [\text{false}] = 0$

- random variable that depends on the sample (just like $\hat{\mu}$)

같은 hypothesis 여도 sample 바라치면
error rate도 달라진다.

Out-of-sample error $E_{\text{out}}(h)$

$$E_{\text{out}}(h) = \mathbb{P}[h(\mathbf{x}) \neq f(\mathbf{x})] \quad (4)$$

► corresponds to μ in the bin model

► unknown but not random (just like μ)

- comparison

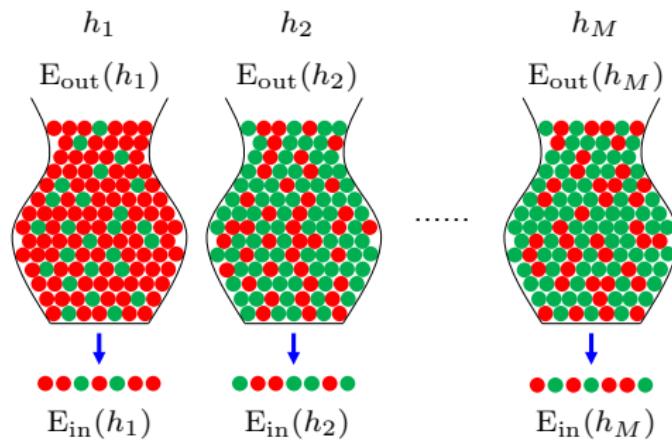
- (E_{in}) is (random) but (known)
- (E_{out}) is (fixed) but (unknown)

Multiple hypotheses

- consider an entire hypothesis set \mathcal{H} instead of just one h
 - ▶ assume for now \mathcal{H} has a *finite* number of hypotheses

$$\mathcal{H} = \{h_1, h_2, \dots, h_M\} \quad M: \text{hypotheses set} \exists \forall \quad (5)$$

- we can then construct a bin equivalent by having M bins



Recall: basic probability rules

for events A and B :

- implication

$$\text{if } A \Rightarrow B \ (A \subseteq B) \text{ then } \mathbb{P}[A] \leq \mathbb{P}[B]$$

- union bound

$$\mathbb{P}[A \text{ or } B] = \mathbb{P}[A \cup B] \leq \mathbb{P}[A] + \mathbb{P}[B]$$

- Bayes' rule

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[B|A] \cdot \mathbb{P}[A]}{\mathbb{P}[B]}$$

Uniform version of Hoeffding inequality

- Hoeffding inequality (1) does not apply to multiple bins
 - ▶ the inequality applies to each bin individually
- a simple but crude fix:
 - ▶ \hat{g} has to be one of the h_m 's (regardless of \mathcal{A} and \mathcal{D})
 - ▶ therefore

$$\begin{aligned} |\text{E}_{\text{in}}(g) - \text{E}_{\text{out}}(g)| > \epsilon &\implies |\text{E}_{\text{in}}(h_1) - \text{E}_{\text{out}}(h_1)| > \epsilon \\ &\text{or } |\text{E}_{\text{in}}(h_2) - \text{E}_{\text{out}}(h_2)| > \epsilon \\ &\dots \\ &\text{or } |\text{E}_{\text{in}}(h_M) - \text{E}_{\text{out}}(h_M)| > \epsilon \end{aligned}$$

- using the union bound

$$\begin{aligned}
 \mathbb{P}[|\text{E}_{\text{in}}(g) - \text{E}_{\text{out}}(g)| > \epsilon] &\leq \mathbb{P}[\quad |\text{E}_{\text{in}}(h_1) - \text{E}_{\text{out}}(h_1)| > \epsilon \\
 &\quad \text{or } |\text{E}_{\text{in}}(h_2) - \text{E}_{\text{out}}(h_2)| > \epsilon \\
 &\quad \dots \\
 &\quad \text{or } |\text{E}_{\text{in}}(h_M) - \text{E}_{\text{out}}(h_M)| > \epsilon] \\
 &\leq \sum_{m=1}^M \mathbb{P}[|\text{E}_{\text{in}}(h_m) - \text{E}_{\text{out}}(h_m)| > \epsilon] \\
 &\leq 2\textcolor{red}{M}e^{-2\epsilon^2 N}
 \end{aligned}$$

- in the end we get a ‘uniform’ version of (2):

$$\mathbb{P}[|\text{E}_{\text{in}}(g) - \text{E}_{\text{out}}(g)| > \epsilon] \leq 2\textcolor{red}{M}e^{-2\epsilon^2 N} \tag{6}$$

Downside for uniform estimates

- probability bound $2Me^{-2\epsilon^2 N}$
 1. factor of M looser than the bound for a single hypothesis
 2. only meaningful if M is finite



- we will improve on these limitations in later lectures

Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound



Feasibility of Learning

러닝이 가능한가?

Error and Noise

Error Measures

Noisy Targets

Summary

Ultimate goal of learning

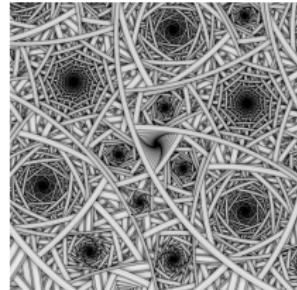
- learning produces a hypothesis g to approximate unknown f
 - ▶ if successful, then g should approximate f well
 - ▶ i.e. $\underline{E_{\text{out}}(g) \approx 0}$ ← the ultimate goal of learning
- however, this is not what we get from the probabilistic analysis
 - ▶ what we get instead: $E_{\text{out}}(g) \approx E_{\text{in}}(g)$
- to conclude that $E_{\text{out}}(g) \approx 0$ ← what we cannot ascertain
 - ▶ we still have to make $\underline{E_{\text{in}}(g) \approx 0}$ (가능) ← what we can ascertain
 - ⇒ we can use E_{in} as a proxy for E_{out}
(대리인)
:裴濬이는 E-test을 대리인으로 사용할 것.

The feasibility of learning

- is split into “two” questions:
 1. can we make sure that $\underline{E_{out}(g) \approx E_{in}(g)}$?
 2. can we make $\underline{E_{in}(g)}$ small enough?
- Hoeffding inequality: $\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$
 - ▶ addresses the first question only
- the second question is answered after we
 - ▶ run the learning algorithm on the actual data, and
 - ▶ see how small we can get E_{in} to be

Complexity of components of learning

- breaking down the feasibility of learning into these two questions
 - ▶ provides further insight into the role of different components
- one such insight is
 - ▶ about the "complexity" of these components
- we consider two types:
 - ▶ complexity of hypothesis set $\mathcal{H} : M$
 - ▶ complexity of target function f



Complexity of hypothesis set \mathcal{H}

- to measure the 'complexity' of \mathcal{H}
 - ▶ we use M the number of hypotheses in \mathcal{H}
- for $E_{\text{out}}(g) \approx E_{\text{in}}(g)$: keep \mathcal{H} less complex
 - ▶ if M goes up \rightarrow we run more risk that $E_{\text{in}}(g)$ will be a poor estimator of $E_{\text{out}}(g)$ [according to (6)]
- for $E_{\text{in}}(g) \approx 0$: better if \mathcal{H} is more complex
 - ▶ g has to come from \mathcal{H} \rightarrow more complex \mathcal{H} gives more flexibility in finding g that fits data well for small $E_{\text{in}}(g)$
- this trade-off in the complexity of \mathcal{H}
 - ▶ a major theme in learning theory

Complexity of target function f

- recall

1. can we make sure that $E_{\text{out}}(g) \approx E_{\text{in}}(g)$?
2. can we make $E_{\text{in}}(g)$ small enough?

- intuitively, a complex f should be harder to learn than a simple f
 - ▶ let's see if this can be inferred from the two questions above
- Hoeffding inequality: $\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$
 - ▶ independent of complexity of f \rightarrow fool independent.
 - ▶ but this doesn't mean: we can learn complex f as easily as simple f
 - ▶ this inequality affects the first question only

- if f gets more complex \rightarrow the second question comes into play
 - ▶ data from more complex f : harder to fit
 \Rightarrow higher $E_{\text{in}}(g)$
- we might try more complex \mathcal{H}
 - ▶ we can fit the data better \Rightarrow lower $E_{\text{in}}(g)$
 - ▶ but then $E_{\text{out}}(g)$ will not be as close to $E_{\text{in}}(g)$ due to (6)
- simply put
 - ▶ a complex f is harder to learn as we expected

Preview: what the theory will allow us to understand

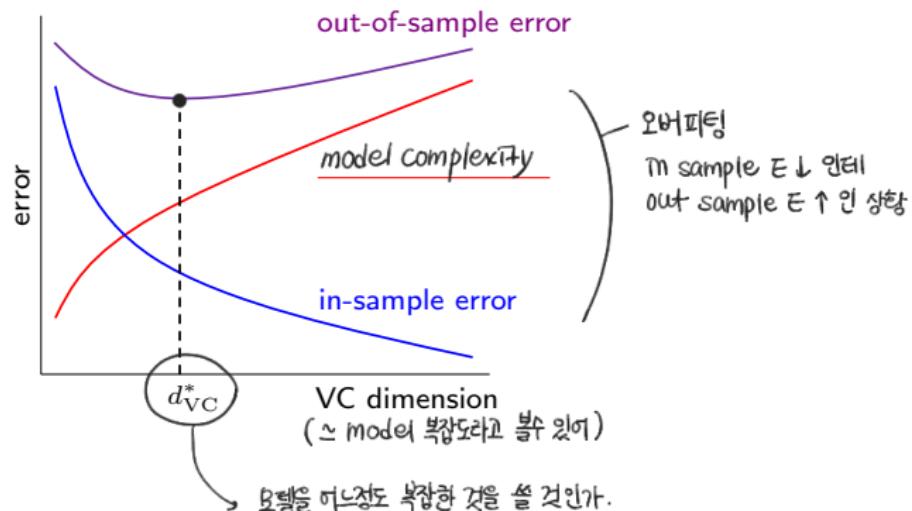
- feasibility of learning for infinite M

무한개 원소 종류가져

(그냥 휘프팅 쓰면 학습 불가능)

- tradeoff:

$(M \uparrow)$	
model complexity \uparrow E_{in}	
model complexity \uparrow $E_{out} - E_{in}$	



Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

Error Measures

Noisy Targets

Summary

Error and noise

- error (aka cost, objective, risk):
 - ▶ about what " $h \approx f$ " means ($g \approx f$)
 - ▶ quantifies how far we are from the target
 - ▶ the choice of an error measure affects the outcome of learning
- noise: (error et 달라)
 - ▶ about the nature of the target function
 - ▶ often noise makes the output of f not uniquely determined by the input
→ fit fn of other conditional probability.



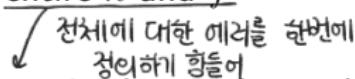
noisy Lena

Pointwise error

- in principle: an error measure quantifies
 - ▶ how well each hypothesis h approximates f

$$\text{Error} = E(h, f) \quad (7)$$

- ▶ this is based on the entire h and f


전체에 대한 에러를 한번에
정의하기 힘들어
- in practice: an error measure is almost always defined
.. 한 점씩 에러를 정의
▶ pointwise on individual input points:

$$e(h(\mathbf{x}), f(\mathbf{x}))$$

- ▶ examples:

squared error $e(h(\mathbf{x}), f(\mathbf{x})) = (h(\mathbf{x}) - f(\mathbf{x}))^2$

binary error $e(h(\mathbf{x}), f(\mathbf{x})) = \llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket$

From pointwise to overall

- overall error $E(h, f) = \underbrace{\text{average}}_{\bigcup} \text{ of pointwise errors } e(h(\mathbf{x}), f(\mathbf{x}))$
- in-sample error:

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N e(h(\mathbf{x}_n), f(\mathbf{x}_n)) \quad (8)$$

- out-of-sample error:

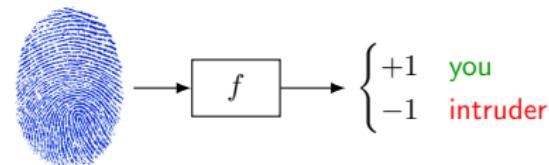
$$E_{out}(h) = \underbrace{\mathbb{E}_{\mathbf{x}}}_{\text{expectation}} [e(h(\mathbf{x}), f(\mathbf{x}))] \quad (9)$$

How to choose the error measure

- ex) fingerprint verification

- two types of error

1. false accept
2. false reject



- how to penalize each type?

		f (target)	
		+1	-1
h (hypothesis)	+1	no error	false accept
	-1	false reject	no error

For supermarkets

- supermarket verifies fingerprint for discounts
- false reject is costly
 - customer gets annoyed!
- false accept is minor
 - gave away a discount and
 - intruder left their fingerprint



		f (target)	
		+1	-1
h (hypothesis)	+1	0	1 (false accept)
	-1	10 (false reject)	0

For the CIA

- CIA verifies fingerprint for security
- false accept is a disaster!
- false reject can be tolerated
 - ▶ try again; you are an employee



		f (target)	
		+1	-1
h (hypothesis)	+1	0	1000 (false accept)
	-1	1 (false reject)	0

Take-home message



(문제에 따라 다르다)

- the error measure should be specified by the user

▶ but this is not always possible

1. the user may not provide an error specification
2. the weighted cost may be a difficult objective function for optimizers to work with

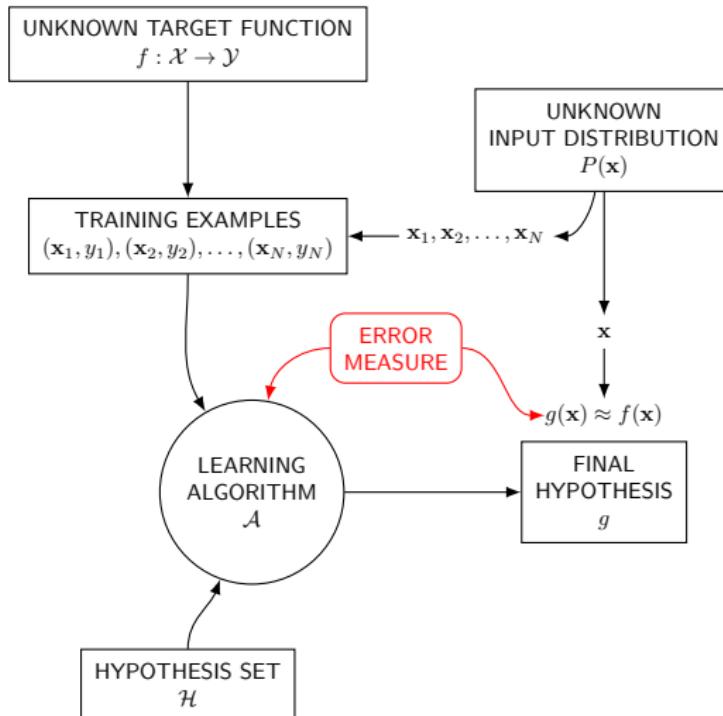
- alternatives

(사람들이 흔히 사용하는 것)

▶ plausible measures: squared error, Gaussian noise

▶ friendly measures: closed-form solution, convex optimization

The learning diagram: with error measures



Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

Error Measures

Noisy Targets

noise: target fn의 뿔

Summary

Noisy targets

- the 'target function' is not always a *function*
- consider the credit-card approval:

feature	value
age	23 years
gender	female
annual salary	\$30,000
years in residence	1 year
years in job	1 year
current debt	\$15,000
...	...

- two 'identical' customers but two different behaviors

- usually, the data we learn from
 - ▶ not generated by a "deterministic target function"
정해지 하지 X : 확률로 생각하자.
 - ▶ instead, generated in a noisy way such that
the output is not uniquely determined by the input
- ex) the credit 'function' distribution
 - ▶ is not really a deterministic function but a noisy one
 - ▶ but still can be modeled within the framework we have

Interpolation versus regression

- given training data $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$

- interpolation**
(not learning)

- $y_n = f(\mathbf{x}_n)$
- find a deterministic function $f(\mathbf{x}) \in \mathbb{R}$ such that $y_n = f(\mathbf{x}_n)$

- regression:**
noise ϵ is added

- $y_n = f(\mathbf{x}_n) + \epsilon$ (한 X 에 대해 다른 y 가 나옵니다)
- find a model $g(\mathbf{x})$ that approximates y_n

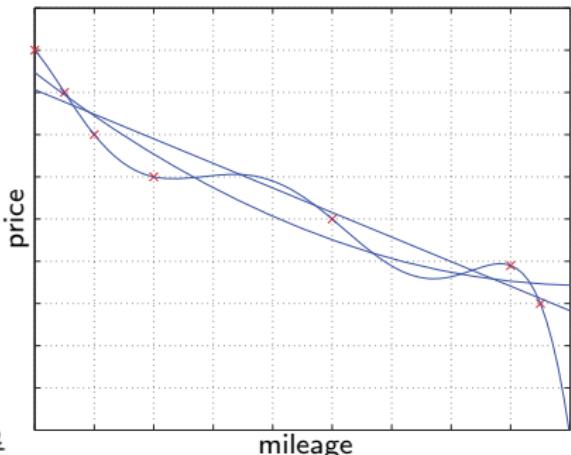


Figure 2: linear, second-order, and sixth-order polynomial fits ($N = 7$)

Target ‘distribution’

- instead of $y = f(\mathbf{x})$, we can take output y to be a random variable
 - ▶ that is affected by, rather than determined by, the input \mathbf{x}
- formally,
 - ▶ we have a target distribution $y \sim p(y|\mathbf{x})$
 - ▶ instead of a target function $y = f(\mathbf{x})$

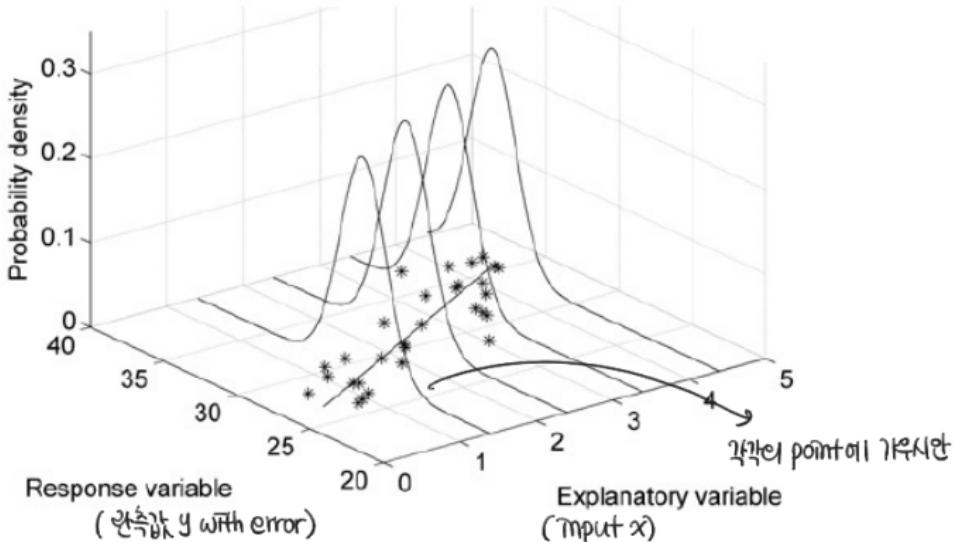


Figure 3: A schematic diagram of a linear regression analysis, where data ($n = 30$, indicated by *) were generated from the linear function with the slope = 2.5, intercept = 25.0, and standard deviation = 1.5. The estimated parameters using a linear regression were slope = 2.45, intercept = 25.3, and standard deviation = 1.48. The estimated regression line is shown as a straight line, whereas four normal distributions indicate the assumed normal distributions of data at four values of the explanatory variable (1.4, 2.4, 3.4, and 4.4).

source: <http://www.seaturtle.org/mtn/archives/mtn122/mtn122p1.shtml>

We can think of

- noisy target as deterministic target plus added noise:

- ▶ noisy target $y = \underbrace{f(\mathbf{x})}_{\text{deterministic target } = \mathbb{E}(y|\mathbf{x})} + \underbrace{y - f(\mathbf{x})}_{\text{pure noise}}$

- deterministic target as the average of noisy target:

- ▶ $f(\mathbf{x}) = \mathbb{E}(y|\mathbf{x})$

- deterministic target as a special case of noisy target:

- ▶ $P(y|\mathbf{x})$ is zero except for $y = f(\mathbf{x})$

The learning diagram: including noisy target

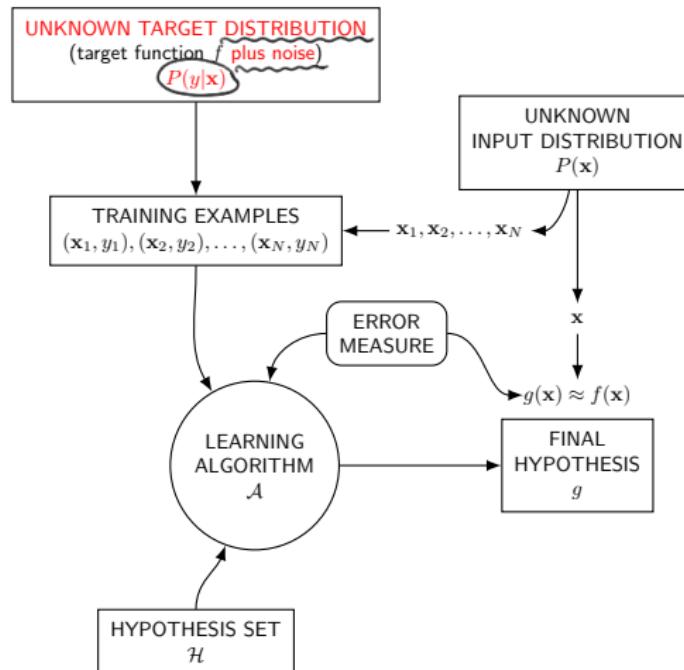


Figure 4: the general (supervised) learning problem

Distinction between $P(y|x)$ and $P(x)$

- both convey probabilistic aspects of x and y
 - ▶ but different roles of $\underbrace{P(y|x)}$ and $\underbrace{P(x)}$ in learning
- $(P(y|x))$, the target distribution
 - ▶ what we are trying to learn
- $(P(x))$, the input distribution
 - ▶ only quantifies the relative Importance of the point x
 $P(x)$ 을 구하는 것도 학습
- merging $P(x)P(y|x)$ as $(P(x,y))$ mixes the two concepts
 - ▶ sample (x, y) is now generated by the joint distribution:
 $\underbrace{P(x)P(y|x)}$
 $P(x,y)$ 은 model이라 한다.

Outline

Learning Unknown Target Function

- Hoeffding Inequality

- Connection to Learning

- Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

- Error Measures

- Noisy Targets

Summary

Summary

- learning *unknown* target f : only possible in a probabilistic sense
 - ▶ we need hypothesis $g \approx f$, which means $\text{E}_{\text{out}}(g) \approx 0$
- feasibility of learning is split into two questions
 1. can we make sure that $\text{E}_{\text{out}}(g) \approx \text{E}_{\text{in}}(g)$?
 2. can we make $\text{E}_{\text{in}}(g) \approx 0$?
- answering Q1: theoretical (first half of this course)
 - ▶ Hoeffding's inequality: $\mathbb{P}[|\text{E}_{\text{in}}(g) - \text{E}_{\text{out}}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$
 - ▶ VC analysis and bias-variance analysis
- answering Q2: practical (second half)
 - ▶ various learning paradigms will help us achieve that

- error measure quantifies the meaning of \approx in expression $g \approx f$
 - ▶ ideally, should be specified by the user
 - ▶ practical alternatives: *plausible/friendly measures*
- averaging *pointwise* error gives
 - ▶ in-sample: $E_{in}(h) = \frac{1}{N} \sum_n e(h(\mathbf{x}_n), f(\mathbf{x}_n))$
 - ▶ out-of-sample: $E_{out}(h) = \mathbb{E}_{\mathbf{x}}[e(h(\mathbf{x}), f(\mathbf{x}))]$
- noisy targets
 - ▶ $y \sim P(y|\mathbf{x})$ instead of $y = f(\mathbf{x})$ [y is a **random variable**]
 - ▶ $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ generated by $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$
 - ▶ $E_{out}(h) = \mathbb{E}_{\mathbf{x}, y}[e(h(\mathbf{x}), y)]$