

Chapter 2: Training versus Testing, Theory of generalization:

1. E_{in} from in sample points; E_{out} from fresh points. E_{in} not always generalize to E_{out} .
Generalization error: discrepancy between E_{in} & E_{out}
2. Error bar depends on infinite M , so is meaningless.
But it is a union bound and overestimated.

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq \frac{2Me^{-2\epsilon^2 N}}{\delta}, \text{ a tolerance level } \delta, \frac{1}{\delta} = 2Me^{-2\epsilon^2 N} \quad E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}.$$

$$E_{out} \geq E_{in} - \epsilon \text{ for all } h \in \mathcal{H}.$$

3. Effective number of hypotheses

➤ Growth function def:

Definition 2.1. Let $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}$. The dichotomies generated by \mathcal{H} on these points are defined by

$$\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \{ (h(\mathbf{x}_1), \dots, h(\mathbf{x}_N)) \mid h \in \mathcal{H} \}. \quad (2.3)$$

Definition 2.2. The growth function is defined for a hypothesis set \mathcal{H} by

$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}} |\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)|,$$

where $|\cdot|$ denotes the cardinality (number of elements) of a set.

pick the N points that give the most dichotomy

Growth function: a measure of number of hypotheses in a set, considered on N points

rather than entire input field. $m_{\mathcal{H}}(N) \leq 2^N$.

➤ Growth function Examples:

Illustrate Example 2.1: growth function for 2D perceptron

Illustrate Example 2.2: cases: positive ray; positive intervals; convex set; hypothesis set complexity \uparrow , growth function \uparrow grows with N faster

➤ Break point definition:

Definition 2.3. If no data set of size k can be shattered by \mathcal{H} , then k is said to be a break point for \mathcal{H} .

If k is a break point, then $m_{\mathcal{H}}(k) < 2^k$. Example 2.1 shows that $k = 4$ is a

Exercise 2.1: find break point of the examples above

[Learning-From-Data-A-Short-Course/Solutions to Chapter 2 Training versus Testing.ipynb at master · josiah Davis/learning-from-data-a-short-course · GitHub](https://github.com/josiah Davis/learning-from-data-a-short-course/blob/master/Chapter%20Training%20versus%20Testing.ipynb)

➤ Bound the growth function:

if no break point, growth func = 2^N , replace M by 2^N in the bound of generalization error (discrepancy between E_{in} and E_{out}), as $N \uparrow$, the bound will not \downarrow (converge to a constant, verifiable)

if growth func bound by a polynomial, replace it in the generalization error, bound of generalization error: numerator: $\ln(N)$; denominator: N . N tends to infinity, bound of error tends to 0. \rightarrow generalize well

Mathematical induction proves the polynomial bound of growth func:

Definition 2.4. $B(N, k)$ is the maximum number of dichotomies on N points such that no subset of size k of the N points can be shattered by these dichotomies.

Lemma 2.3 (Sauer's Lemma).

$$B(N, k) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

Theorem 2.4. If $m_{\mathcal{H}}(k) < 2^k$ for some value k , then

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

for all N . The RHS is polynomial in N of degree $k - 1$.

Homework: exercise 2.2 verification of theorem 2.4

$$\sum_0^n \binom{n}{i} = 2^n \quad (1+1)^n = \sum_{i=0}^n \binom{n}{i} 1^i 1^{n-i}.$$

Handwritten derivation showing that $m_{\mathcal{H}}(N) = N + 2^{\lfloor \frac{k}{2} \rfloor} \leq \sum_{i=0}^{k-1} \binom{N}{i}$; k : break point. It then shows the sum is a polynomial in N of degree $k-1$, which is not true for all N . A note on the right says "or $k = N+1$ (Valid, max)".

4. The VC dimension:

Definition 2.5. The Vapnik-Chervonenkis dimension of a hypothesis set \mathcal{H} , denoted by $d_{\text{VC}}(\mathcal{H})$ or simply d_{VC} , is the largest value of N for which $m_{\mathcal{H}}(N) = 2^N$. If $m_{\mathcal{H}}(N) = 2^N$ for all N , then $d_{\text{VC}}(\mathcal{H}) = \infty$.

$k = d_{\text{VC}} + 1$ \mathcal{H} can shatter d_{VC} points,

Exercise 2.3: $d = k - 1$

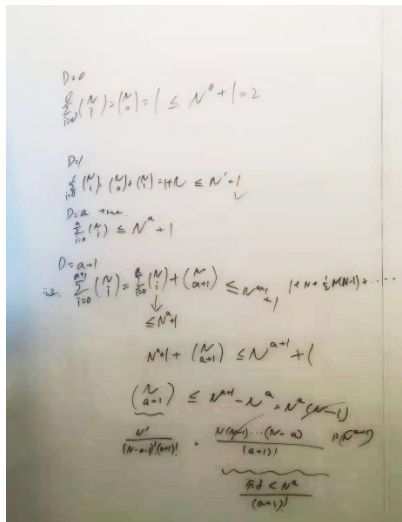
$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{d_{\text{VC}}} \binom{N}{i}.$$

the VC dimension is the order of the polynomial bound on growth func.

$$m_{\mathcal{H}}(N) \leq N^{d_{\text{VC}}} + 1.$$

Problem 2.5 Prove by induction that $\sum_{i=0}^D \binom{N}{i} \leq N^D + 1$, hence

$$m_{\mathcal{H}}(N) \leq N^{d_{\text{VC}}} + 1.$$



或者用 github 的证明

5. New inequality -> VC Bound

$m_H(2N)$:

Hoeffding inequality: entire input space

Growth function: only data set

\therefore replace E_{out} with E_{in} (depend on another N)

\therefore total # of hypotheses depends on $2N$.

Not quite:

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

but rather:

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 4 m_H(2N) e^{-\frac{1}{8}\epsilon^2 N}$$

The Vapnik-Chervonenkis Inequality

Theorem 2.5 (VC generalization bound). For any tolerance $\delta > 0$,

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{N} \ln \frac{4m_H(2N)}{\delta}} \quad (2.12)$$

with probability $\geq 1 - \delta$.

$$\begin{aligned}
P[|E_{in}(g) - E_{out}(g)| > \epsilon] &\leq 4 m_H(2N) e^{-\frac{1}{8} \epsilon^2 N} \\
&\Downarrow \delta \\
E_{out}(g) - E_{in} &\leq \epsilon \quad \text{with probability } 1 - \delta \\
E_{out}(g) &\leq E_{in}(g) + \epsilon \\
\delta &= 4 m_H(2N) e^{-\frac{1}{8} \epsilon^2 N} \\
\frac{\delta}{4 m_H(2N)} &= e^{-\frac{1}{8} \epsilon^2 N} \\
-\frac{1}{8} \epsilon^2 N &= \ln \frac{\delta}{4 m_H(2N)} \\
\epsilon^2 &= \frac{-8 \ln \frac{\delta}{4 m_H(2N)}}{N} = \frac{8}{N} \ln \frac{4 m_H(2N)}{\delta} \\
\epsilon &= \sqrt{\frac{8}{N} \ln \frac{4 m_H(2N)}{\delta}} \\
E_{out}(g) &\leq E_{in}(g) + \sqrt{\frac{8}{N} \ln \frac{4 m_H(2N)}{\delta}}
\end{aligned}$$

6. Why VC bound so loose:

Exercise 2.5: VC bound very loose

The slack is from:

- Hoeffding inequality: different E_{out} \rightarrow different variance in E_{in} , \rightarrow different probability, however, all the cases covered by one bound \therefore hoeffding inequality loose
- Growth function $m_H(N)$: did NOT consider probability distribution on input space (if consider, pick specific data set or use expectation). $m_H(N)$ is an upper bound
- Use the bound of $m_H(N)$ (polynomial in order d_{vc}), rather than its value, further slack

VC analysis (bound & dimension) \rightarrow generalization performance

Both useful in practice

7. Sample complexity:

how many training examples N are needed to achieve a certain generalization performance

error tolerance: ϵ the allowed generalization error

confidence parameter: δ how often the error tolerance ϵ is violated

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{N} \ln \frac{4 m_H(2N)}{\delta}}$$

with probability $\geq 1 - \delta$.

$$\sqrt{\frac{8}{N} \ln \frac{4 m_H(2N)}{\delta}} \leq \epsilon \quad N \geq \frac{8}{\epsilon^2} \ln \left(\frac{4 m_H(2N)}{\delta} \right) \quad N \geq \frac{8}{\epsilon^2} \ln \left(\frac{4((2N)^{d_{vc}} + 1)}{\delta} \right)$$

Get N in an iterative way

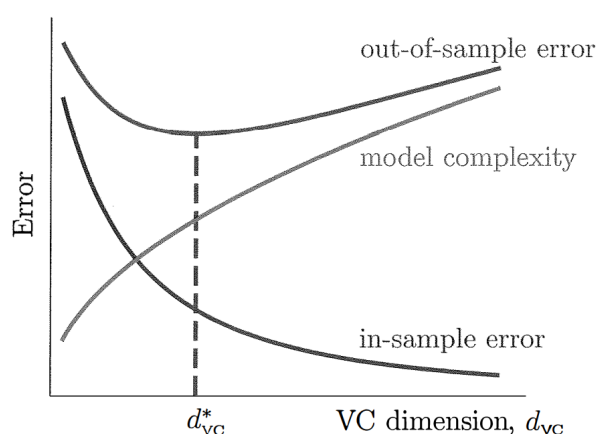
8. penalty for model complexity

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \Omega(N, \mathcal{H}, \delta),$$

where

$$\begin{aligned} \Omega(N, \mathcal{H}, \delta) &= \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)} \\ &\leq \sqrt{\frac{8}{N} \ln \left(\frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} \right)}. \end{aligned}$$

- a) **Higher confidence** (更加保证偏离不超过某个值, error 大于某个值的可能性非常小的 (Hoeffding inequality)), **lower δ** , 这个 error, 或者说 **penalty $\Omega \uparrow$** 。因为要更 confident, 只能放宽条件
- b) More training examples, **$N \uparrow$** , fit better, **penalty $\Omega \downarrow$**
- c) More complex hypothesis set, higher VC dimension, more choices to fit the data **$E_{\text{in}} \downarrow$** , **penalty $\Omega \uparrow$**



9. Etest estimates Eout

Estimate Eout using a test set which is not involved in the training process

Etest generalizes to Eout by Hoeffding inequality since only one hypothesis g applies to Etest
Exercise 2.6

- (a) M is given directly, so do not need to bother with growth func or dvc (applies to binary target function), just use the initial Hoeffding multiply with # of hypothesis
- (b) more examples on training, more complex hypothesis set, may result in greater penalty
Less examples on training, fewer choices to fit the data, maybe less good g , even generalization (to Eout) bound on Etest is small (Eout Etest close), Etest may be very large due to bad g .

[Learning-From-Data-A-Short-Course/Solutions to Chapter 2 Training versus Testing.ipynb at master · niuers/Learning-From-Data-A-Short-Course \(github.com\)](#)

Training set: bias in estimate of Eout, VC bound takes the bias

Test set: no bias, only variance due to finite sample size, just tell how well we did

10. Other Target Types

Real-valued function \rightarrow square error

$$E_{\text{out}}(h) = \mathbb{E}[(h(\mathbf{x}) - f(\mathbf{x}))^2] \quad E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - f(\mathbf{x}_n))^2$$

Exercise 2.7 hint: $\mathbb{E}[\mathbf{x}] = \sum \mathbf{x} P(\mathbf{x})$

Law of Large Number: E_{in} converge to E_{out} . Hoeffding inequality is one form

11. Approximation – generalization tradeoff

a) E_{out} decomposition: bias and variance

$$\begin{aligned} E_{\text{out}}(g^{(\mathcal{D})}) &= \mathbb{E}_{\mathbf{x}} [(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2], \\ \mathbb{E}_{\mathcal{D}} [E_{\text{out}}(g^{(\mathcal{D})})] &= \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\mathbf{x}} [(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2]] \\ &= \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} [(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2]] \\ &= \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2] - 2 \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})] f(\mathbf{x}) + f(\mathbf{x})^2] \end{aligned}$$

Handwritten derivation for $\mathbb{E}_{\mathcal{D}}(g^{(\mathcal{D})}(\mathbf{x}))$:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}(g^{(\mathcal{D})}(\mathbf{x})) &: p_1, \dots, p_k \\ &\downarrow \\ &g_1, \dots, g_k \\ &\downarrow \\ &\frac{1}{k} \sum g_1(\mathbf{x}) \dots g_k(\mathbf{x}) \Rightarrow \bar{g}(\mathbf{x}) \end{aligned}$$

random variable

$$\begin{aligned} &\mathbb{E}_{\mathcal{D}} [E_{\text{out}}(g^{(\mathcal{D})})] \\ &= \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2] - 2 \bar{g}(\mathbf{x}) f(\mathbf{x}) + f(\mathbf{x})^2], \\ &= \mathbb{E}_{\mathbf{x}} \left[\underbrace{\mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2] - \bar{g}(\mathbf{x})^2}_{\textcircled{1}} + \underbrace{\bar{g}(\mathbf{x})^2 - 2 \bar{g}(\mathbf{x}) f(\mathbf{x}) + f(\mathbf{x})^2}_{\textcircled{2}} \right], \end{aligned}$$

$\textcircled{2} \rightarrow \text{bias}$

$$\begin{aligned} \textcircled{1} &= \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2 + \bar{g}(\mathbf{x})^2 - 2 g^{(\mathcal{D})}(\mathbf{x}) \bar{g}(\mathbf{x})] \\ &= \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2] + \bar{g}(\mathbf{x})^2 - 2 \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})] \bar{g}(\mathbf{x}) \\ &= \underbrace{\mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2]}_{\textcircled{1}} + \underbrace{\bar{g}(\mathbf{x})^2}_{\textcircled{2}} - 2 \bar{g}(\mathbf{x})^2 \\ &= \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(\mathbf{x})^2] - \bar{g}(\mathbf{x})^2 \end{aligned}$$

$$\text{bias}(\mathbf{x}) = (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2$$

How much the learning model is biased away from target function
 Limitation from learning model

$$\text{var}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]$$

Variation in the final hypothesis, depending on data set

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[E_{\text{out}}(g^{(\mathcal{D})})] &= \mathbb{E}_{\mathbf{x}}[\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})] \\ &= \text{bias} + \text{var},\end{aligned}$$

Variance: small variation in data \rightarrow vastly different hypothesis

Small model: large bias, small var

Large model: large var, small bias

If data noisy: additional noise term

b) Problem 2.22

Problem 2.22

$$E_{\text{out}}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}, y}[(g^{(\mathcal{D})}(\mathbf{x}) - y(\mathbf{x}))^2] \quad y(\mathbf{x}) = f(\mathbf{x}) + \varepsilon$$

ε zero mean random variable. $\mathbb{E}_{\mathbf{x}}(\varepsilon) = 0$ ~~Var~~ $\text{Var}(\varepsilon) = \sigma^2$

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[E_{\text{out}}(g^{(\mathcal{D})})] &= \mathbb{E}_{\mathcal{D}}[\mathbb{E}_{\mathbf{x}, y}[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) - \varepsilon)^2]] \\ &= \mathbb{E}_{\mathbf{x}}[\underbrace{\mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2]}_{\text{var}(\mathbf{x}) + \text{bias}(\mathbf{x})} + \mathbb{E}_{\mathcal{D}}[\varepsilon^2] - 2\mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))\varepsilon]] \\ &= \mathbb{E}_{\mathbf{x}}[\text{var}(\mathbf{x}) + \text{bias}(\mathbf{x})] + \mathbb{E}_{\mathbf{x}}(\sigma^2) - 2\underbrace{\mathbb{E}_{\mathbf{x}}(\varepsilon)}_{=0} \underbrace{\mathbb{E}_{\mathbf{x}}[g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})]}_{\text{independent.}} \\ &= \text{var} + \text{bias} + \sigma^2\end{aligned}$$

c) Example 2.8

Hypothesis 0: bias & var 的推导

Hypothesis 1: 用代码算

$$f(x) = \sin(\pi x)$$

$$\bar{g}(x) = E_D(g^{(D)}(x))$$

$$\text{hypothesis} = g(x) = \frac{y_1 + y_2}{2}$$

$$y_i = \sin(\pi x_i)$$

$$= E_D\left(\frac{y_1 + y_2}{2}\right)$$

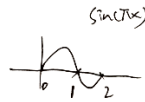
$$= \frac{1}{2} E_D(y_1) + \frac{1}{2} E_D(y_2)$$

$$= \frac{1}{2} E_D(\sin(\pi x_1)) + \frac{1}{2} E_D(\sin(\pi x_2))$$

$$= 0 + 0 = 0$$

$$\text{bias}(x) = (\bar{g}(x) - f(x))^2$$

$$= f(x)^2 = [\sin(\pi x)]^2$$



$$\text{bias} = E_x(\text{bias}(x)) = E_x[\sin^2(\pi x)]$$

$$= \int_0^2 \sin^2(\pi x) \cdot \frac{1}{2} dx$$

$$= \frac{1}{2} \int_0^2 \frac{1 - \cos 2\pi x}{2} dx$$

$$= \frac{1}{2} \int_0^2 \frac{1}{2} dx - \frac{1}{2} \int_0^2 \frac{\cos 2\pi x}{2} dx$$

$$= \frac{1}{2} - \frac{1}{4} \frac{1}{2\pi} (\sin 2\pi x)_0^2$$

$$= \frac{1}{2} - \frac{1}{4} \frac{1}{2\pi} (0 - 0) = \frac{1}{2}$$

$$\cos 2x = 1 - 2\sin^2 x$$

$$\sin^2 x = \frac{1 - \cos 2x}{2}$$

$$\frac{1}{x} \Big|_{-1}^1 = \frac{1}{2} (1 + 1)$$

$$\text{var}(x) = E_D[(g^{(D)}(x) - \bar{g}(x))^2]$$

$$= E_D[g^{(D)}(x)^2]$$

$$= E_D\left[\left(\frac{y_1 + y_2}{2}\right)^2\right]$$

$$= E_D\left[\frac{y_1^2 + y_2^2 + 2y_1 y_2}{4}\right]$$

$$= \frac{1}{4} E_D[y_1^2 + y_2^2 + 2y_1 y_2]$$

$$= \frac{1}{4} \left(\frac{1}{2} + \frac{1}{2}\right) = \frac{1}{4}$$

$$\text{var} = E_x(\text{var}(x)) = 0.25$$

$$\begin{aligned} x_1, x_2 \\ \text{independent} \\ \bar{g}(x) &= 0 \\ y_i &= \sin(\pi x_i) \\ E_D(\sin^2(\pi x_i)) &= \frac{1}{2} \\ E_D(\sin(\pi x_1) \sin(\pi x_2)) \\ &= 0 \times 0 = 0 \end{aligned}$$

For the Visionaries

d) Recap VC analysis:

Growth function: a measure of number of hypotheses in a set, considered on N points
the VC dimension is the **order** of the polynomial bound on growth func.

$$m_{\mathcal{H}}(N) \leq N^{d_{\text{vc}}} + 1.$$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \Omega(N, \mathcal{H}, \delta),$$

where

$$\begin{aligned} \Omega(N, \mathcal{H}, \delta) &= \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)} \\ &\leq \sqrt{\frac{8}{N} \ln \left(\frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} \right)}. \end{aligned}$$

More complex hypothesis set, higher VC dimension, more choices to fit the data $E_{\text{in}} \downarrow$, penalty $\Omega \uparrow$

(Growth function & VC Dimension): Given N points (a set), how many hypothesis \rightarrow complexity of hypothesis

- e) Example 2.8: given 2 points, the two models (with different learning algorithm), both have one possible hypothesis each (complexity of hypothesis same) \rightarrow same VC bound on out-of-sample error.
- f) VC analysis: depend purely on Hypothesis Set, independent of learning algorithm
Bias-Variance analysis: out of sample error: bias and variance decomposition

$$E_{\text{out}}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right],$$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[E_{\text{out}}(g^{(\mathcal{D})})] &= \mathbb{E}_{\mathbf{x}}[\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})] \\ &= \text{bias} + \text{var}, \end{aligned}$$

$$\text{bias}(\mathbf{x}) = (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2 \quad \text{var}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]$$

$$\bar{g}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[g^{(\mathcal{D})}(\mathbf{x})]$$

Different complexity of hypothesis set / learning algorithm \rightarrow different $g^{(\mathcal{D})}$, building block of bias-variance analysis \rightarrow different **bias** & **var** terms

(Bias-Variance analysis based on squared-error (to measure bias and variance), the learning algorithm not have to base on minimizing squared-error measure)

- g) Application of Bias-variance analysis:

to compute: Target function + input probability distribution \therefore just a conceptual tool to DEVELOP A MODEL

\downarrow bias, -var: prior info regarding the target, to steer the Hypothesis Set in the direction of target func

\downarrow var, -bias: general techniques

h) The learning curve of a model

In Sample Error: $\mathbb{E}_{\mathcal{D}}[\hat{E}_{\text{in}}(g^{(\mathcal{D})})]$

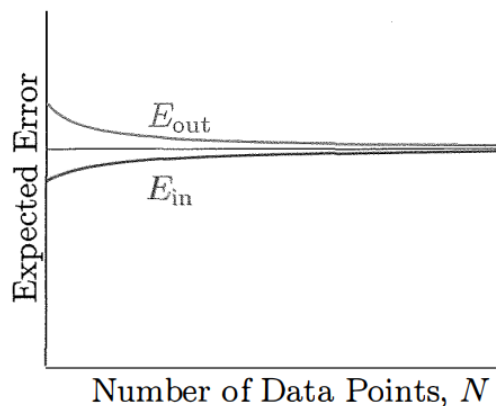
Out of Sample Error: $\mathbb{E}_{\mathcal{D}}[E_{\text{out}}(g^{(\mathcal{D})})]$.

Common:

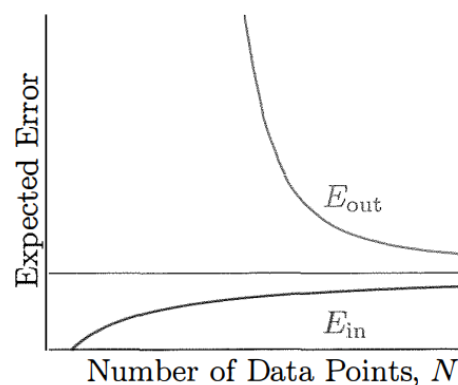
Small N: learning model has a easier task, regardless of outside N. Ein small, Eout large

$N \uparrow$: hard to fit, $E_{\text{in}} \uparrow$ -> smallest error the learning model can achieve (点差不多都出来了)

$N \uparrow$: $E_{\text{out}} \downarrow$ to the best the learning model can achieve



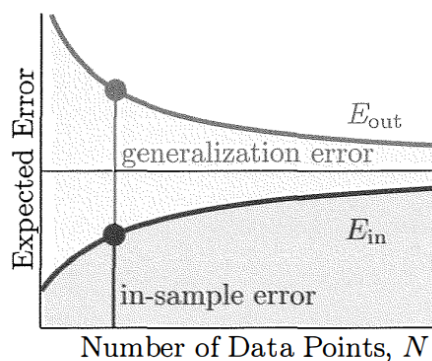
Simple Model



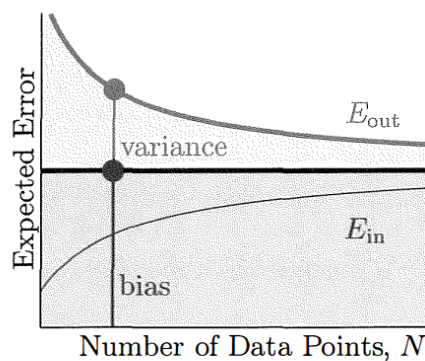
Complex Model

Differences:

Simple model: Converge more quickly but worse ultimate performance



VC Analysis



Bias-Variance Analysis

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \Omega(N, \mathcal{H}, \delta),$$

$$\begin{aligned} \Omega(N, \mathcal{H}, \delta) &= \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)} \\ &\leq \sqrt{\frac{8}{N} \ln \left(\frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} \right)}. \end{aligned}$$

$$E_{\text{out}} = \text{bias} + \text{var}$$

Assumption: best approximation to f of the model

$$= \text{bias} = \bar{g}(\mathbf{x}) \quad \text{performance}$$

Generalization error bounded by Ω (penalty for model complexity)

$N \uparrow, \Omega \downarrow$, generalization error \downarrow

Recap:

1. Given the Hoeffding inequality, how can it be modified to apply a bound to a hypothesis set with M terms?

2. What is the bound of E_{out} with respect to E_{in} , suppose the confidence that E_{out} is within the bound is $1 - \delta$

3. Replace the bound you derived above by VC generalisation bound.

What does each term represent?

growth function - formalise the effective number of hypothesis

4. sample size, complexity of hypothesis set, tolerance level δ .

How are the parameters influencing the VC generalisation bound?

5. How is the VC dimension (model complexity) affecting the E_{in} and E_{out} , draw a graph

6. Based on square error measures, what is the out-of-sample error?

(hint: make explicit the dependence of the final hypothesis on its particular data set used)

7. This out-of-sample error can be decomposed into "bias" and "variance", what do bias and variance imply here?

8. Plot and illustrate the "Learning Curves" of simple model and complex model respectively.

9. Illustrate the VC analysis and Bias-Variance analysis with the learning curve