# BIAS-VARIANCE TRADEOFF

Sonpvh

# OUTLIER

# 1. LEARNING FROM DATA [2]



UNKNOWN TARGET FUNCTION
$f: X \rightarrow Y$
*(ideal credit approval function)*

TRAINING EXAMPLES
$(x_1, y_1), \dots, (x_N, y_N)$
*(historical records of credit customers)*

LEARNING ALGORITHM
$\mathcal{A}$

FINAL HYPOTHESIS
$g \approx f$
*(final credit approval formula)*

HYPOTHESIS SET
$\mathcal{H}$
*(set of candidate formulas)*

1. **Learning:**

   1. Unknown target function $y = f(x)$

   2. Dataset $(x_1, y_1)$, $(x_2, y_2)$ …

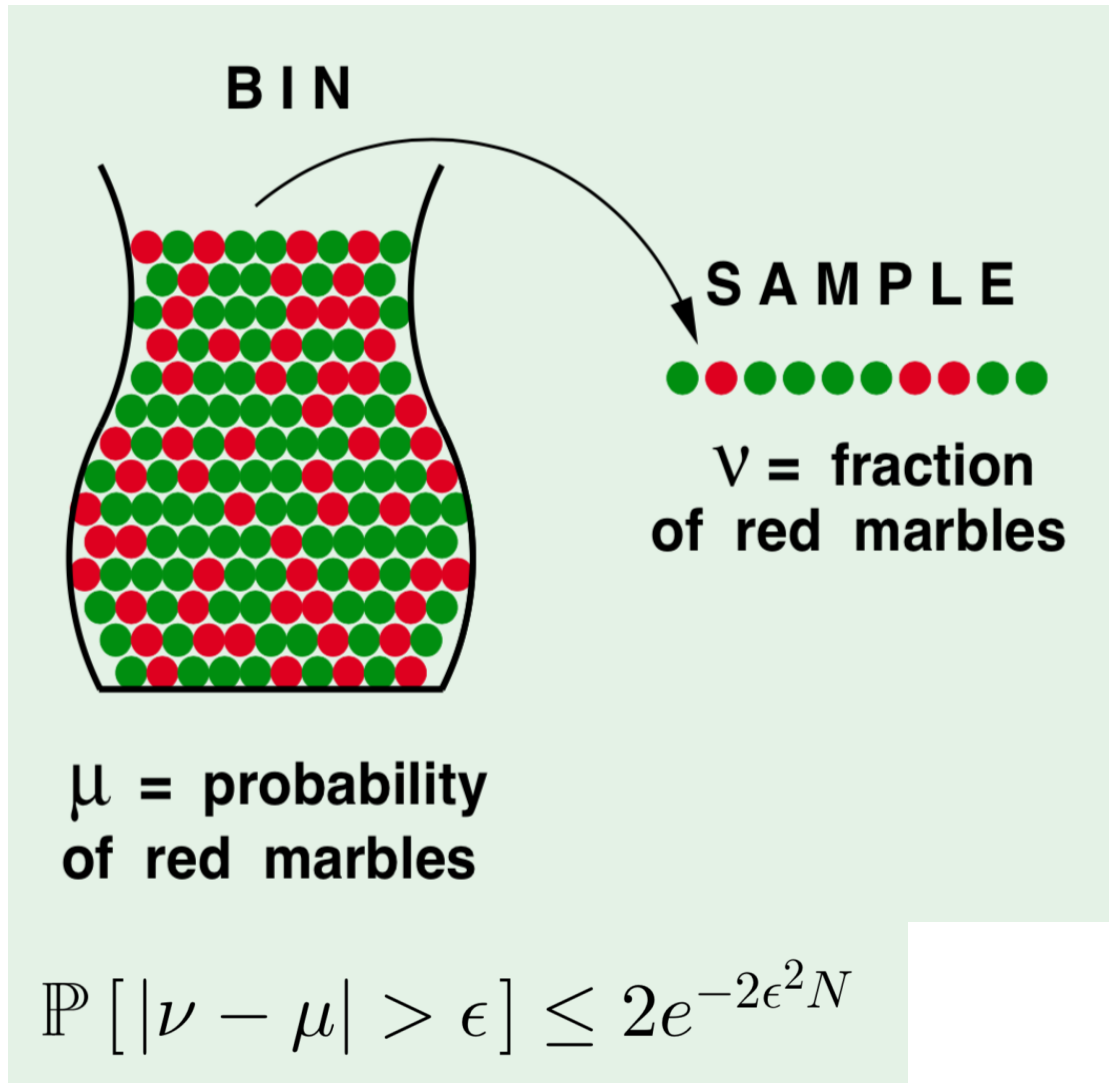   3. Learning algorithms pick $g \sim f$ from Hypothesis Set H

2. **Learning components:**
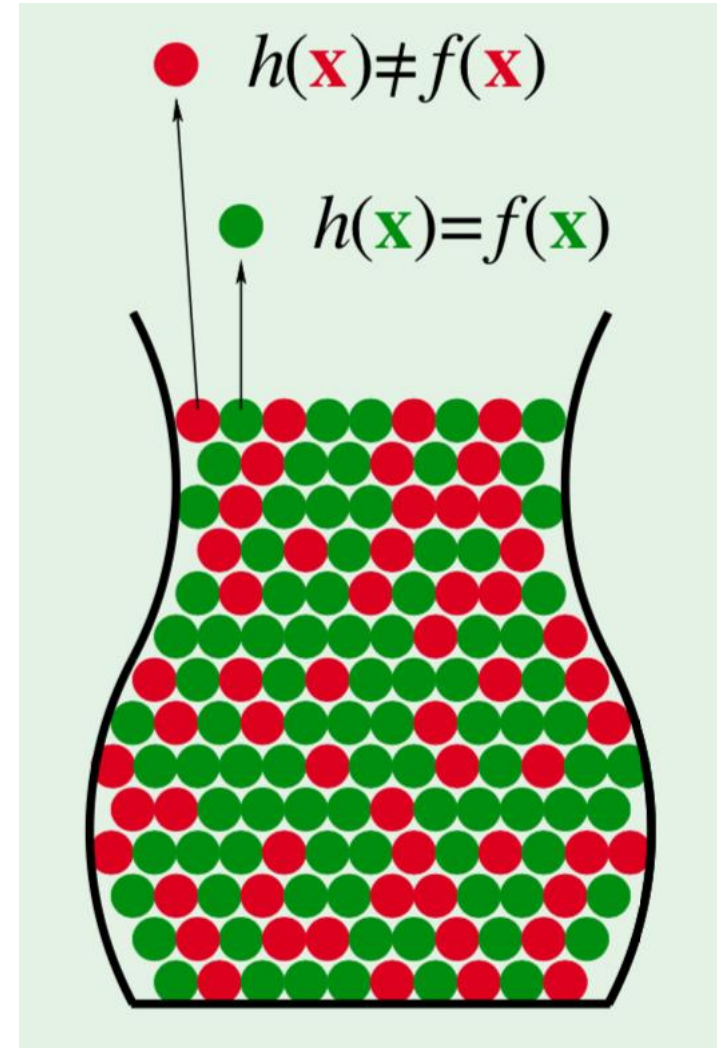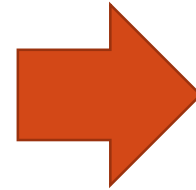
   1. Learning algorithm

   2. Hypothesis set

3. **Purpose:**

   1. $g(x) \sim f(x)$

2

# 1. IS LEARNING FEASIBLE [2]



**BIN**

**SAMPLE**

$\nu$ = fraction of red marbles

$\mu$ = probability of red marbles

$$\mathbb{P}\left[|\nu - \mu| > \epsilon\right] \le 2e^{-2\epsilon^2 N}$$
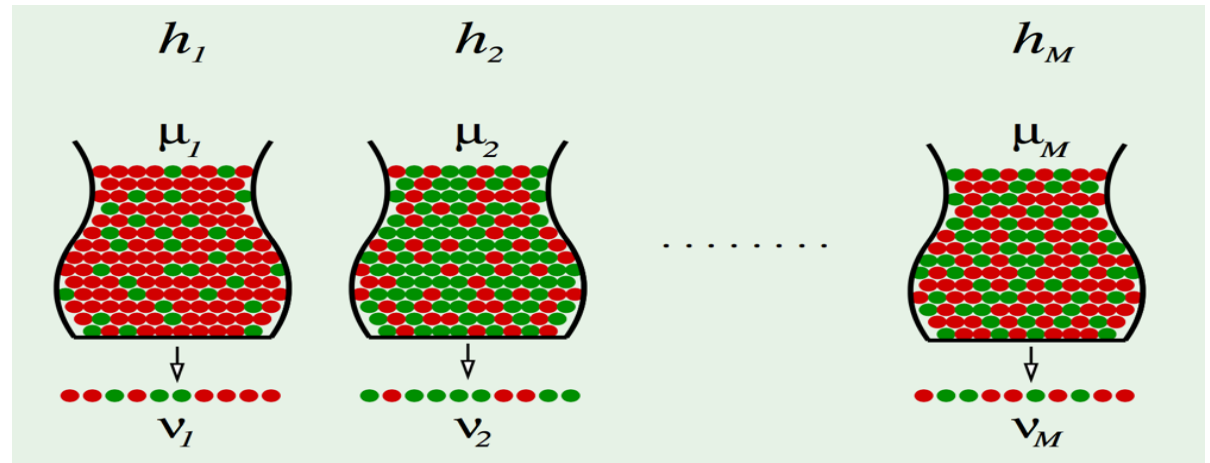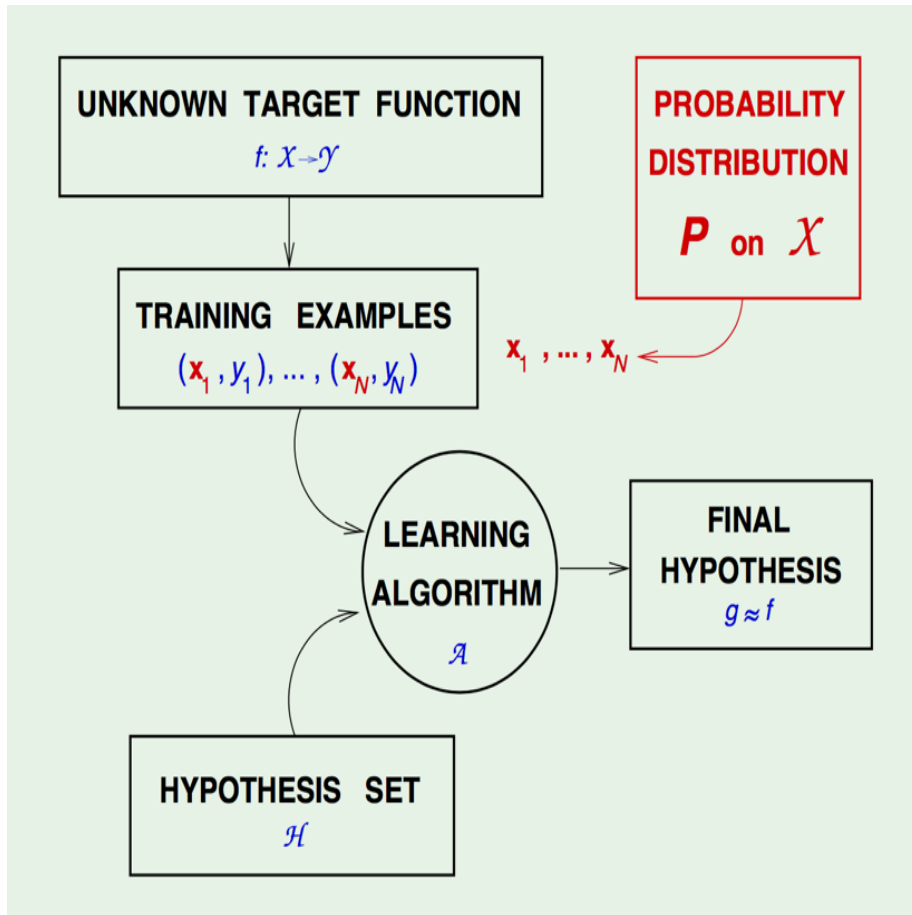
$h(\mathbf{x}) \not\equiv f(\mathbf{x})$

$h(\mathbf{x}) = f(\mathbf{x})$
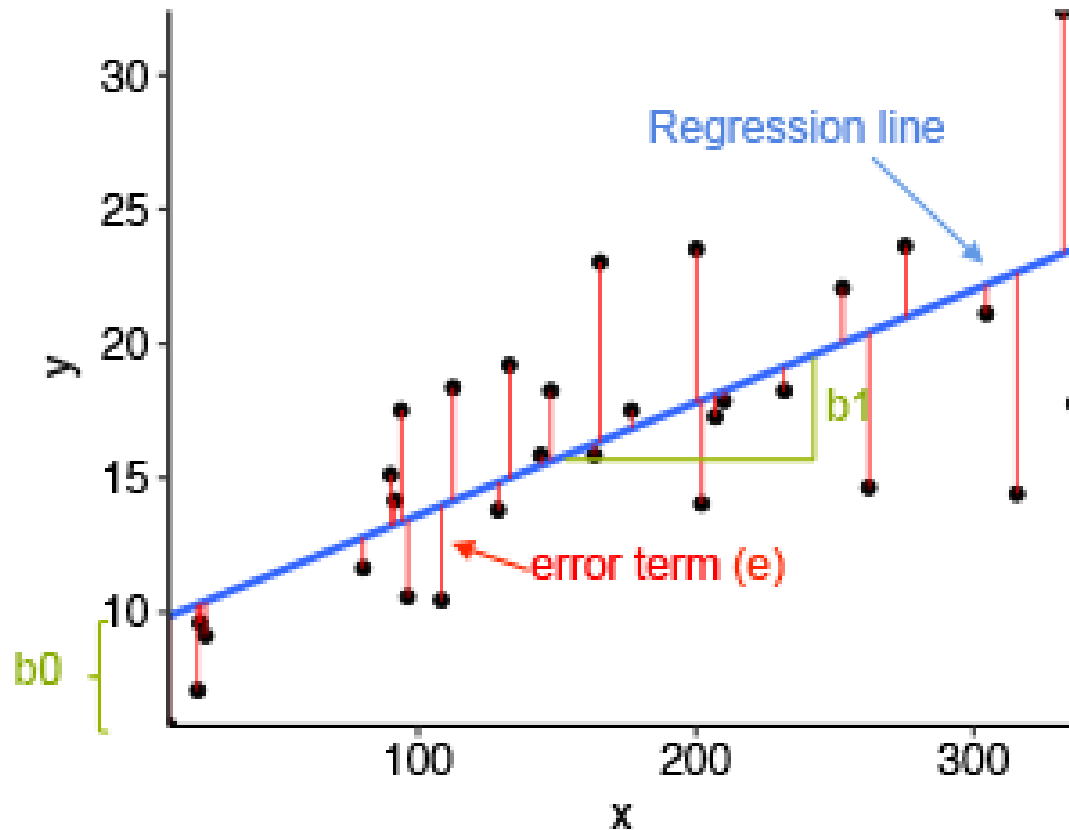
Hoeffding's inequality

# 1. IS LEARNING FEASIBLE [2]



$$\mathbb{P}\big[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\big] \le 2Me^{-2\epsilon^2 N}$$

# 2. ERROR [2]



**Learning Purpose:** $g(x) \sim f(x)$

But what the "$g \sim f$" mean ?  $E(g,f)$

Squared error:      $e\left(h(\mathbf{x}), f(\mathbf{x})\right) = \left(h(\mathbf{x}) - f(\mathbf{x})\right)^2$

Binary error:        $e\left(h(\mathbf{x}), f(\mathbf{x})\right) = [\![h(\mathbf{x}) \neq f(\mathbf{x})]\!]$

[1]

Supper market
verify for discount

|     |     | $f$ | |
| --- | --- | --- | --- |
|     |     | $+1$ | $-1$ |
| $h$ | $+1$ | $0$ | $1$ |
|     | $-1$ | $10$ | $0$ |

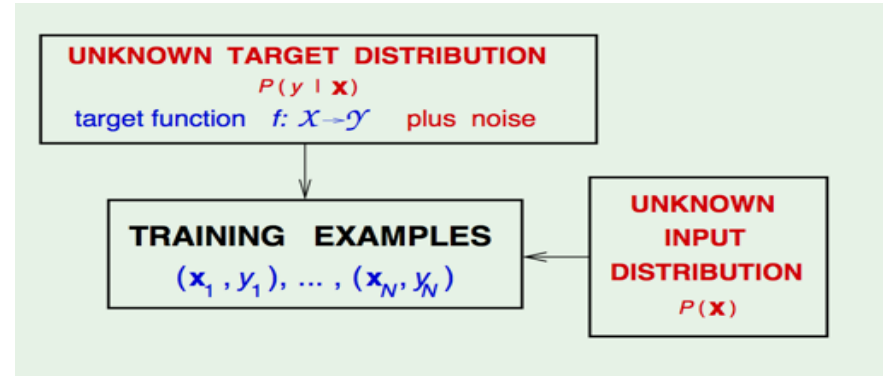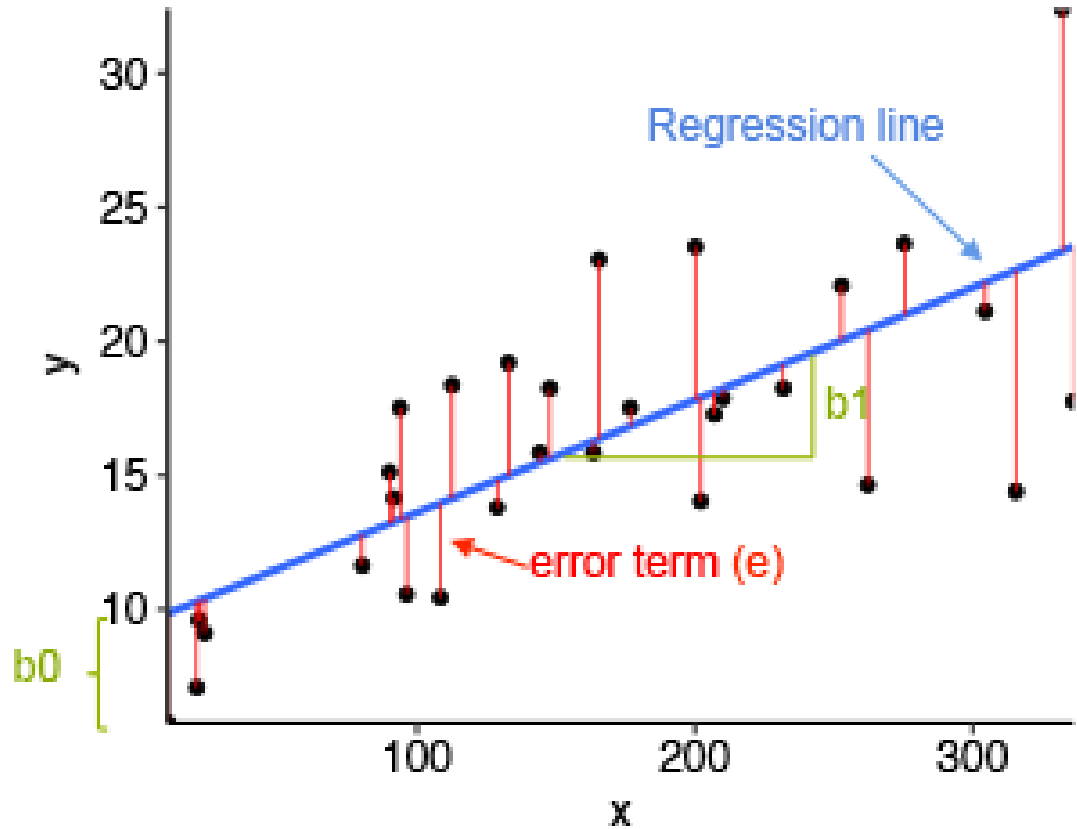|     |     | $f$ | |
| --- | --- | --- | --- |
|     |     | $+1$ | $-1$ |
| $h$ | $+1$ | $0$ | $1000$ |
|     | $-1$ | $1$ | $0$ |

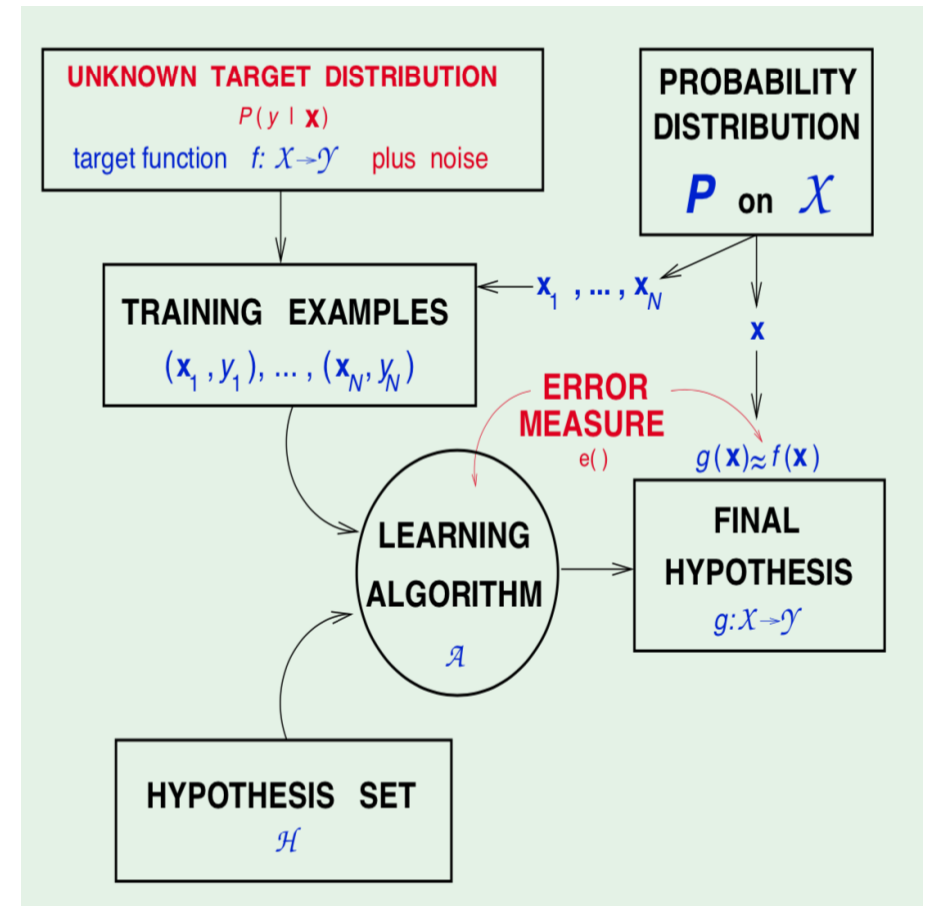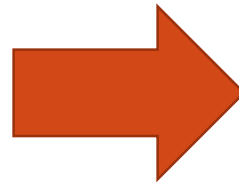CIA verify for security
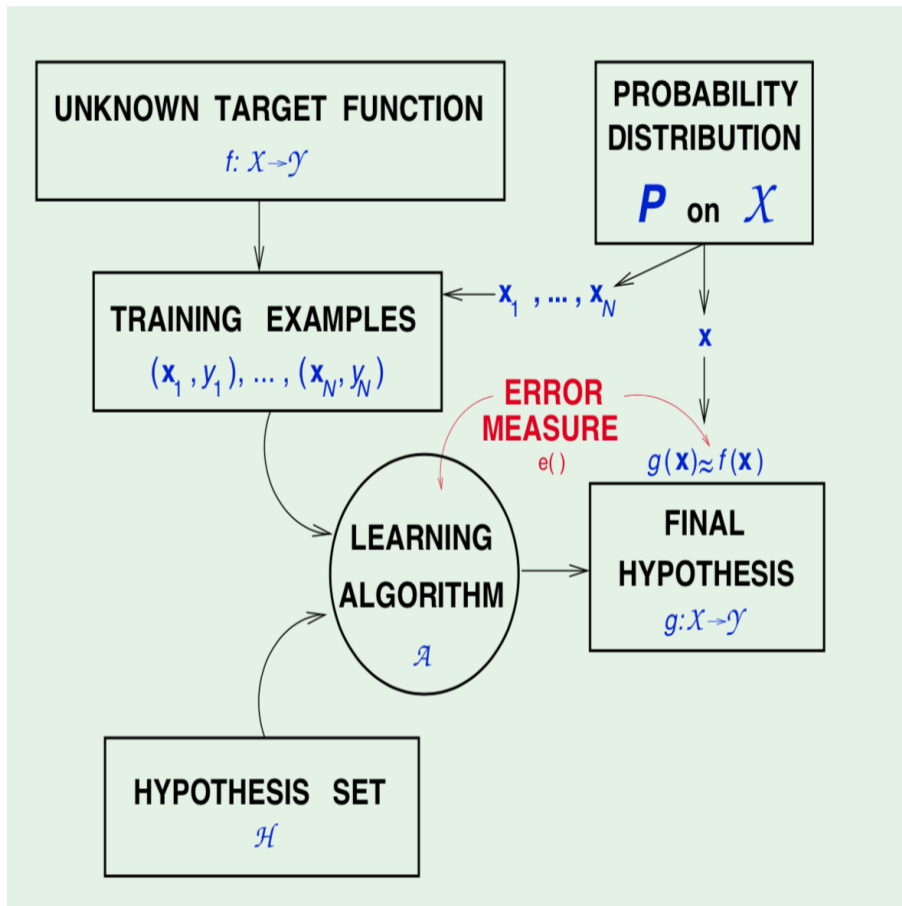
# 2. ERROR [2]

# 2. NOISE [2]





$$y = \hat{y} + \text{noise} = f(x) + \text{noise} = \mathbb{E}\,(y|x) + \text{noise}$$

# 2. NOISE [2]

# 2. PREAMBLE OF THE THEORY [2]

$$E_{out}(g) \approx E_{in}(g) \ (1)$$

$$E_{in}(g) \approx 0 \ (2)$$

(1) Hoeffding's inequality

(2) Optimize error

➔ g ~ f

- f(x) – y = (stochastic) noise
- f(x) – g(x) = (deterministic) noise
- y – g(x) = error

# 3. APPROXIMATION-GENERALIZATION [2]

income

expenditure

$E_{in}$

$E_{out}$

$$E_{in}(g) \approx 0$$

$$E_{out}(g) \approx E_{in}(g)$$

Model Complexity $\sim \mathcal{H}$

# 3. APPROXIMATION-GENERALIZATION [2]



$$d_{VC} \sim \mathcal{H}$$

**Approximation – generalization trade-off**

More complex H ➜ better chance of approximation f

Less complex H ➜ better chance of generalizing out of sample
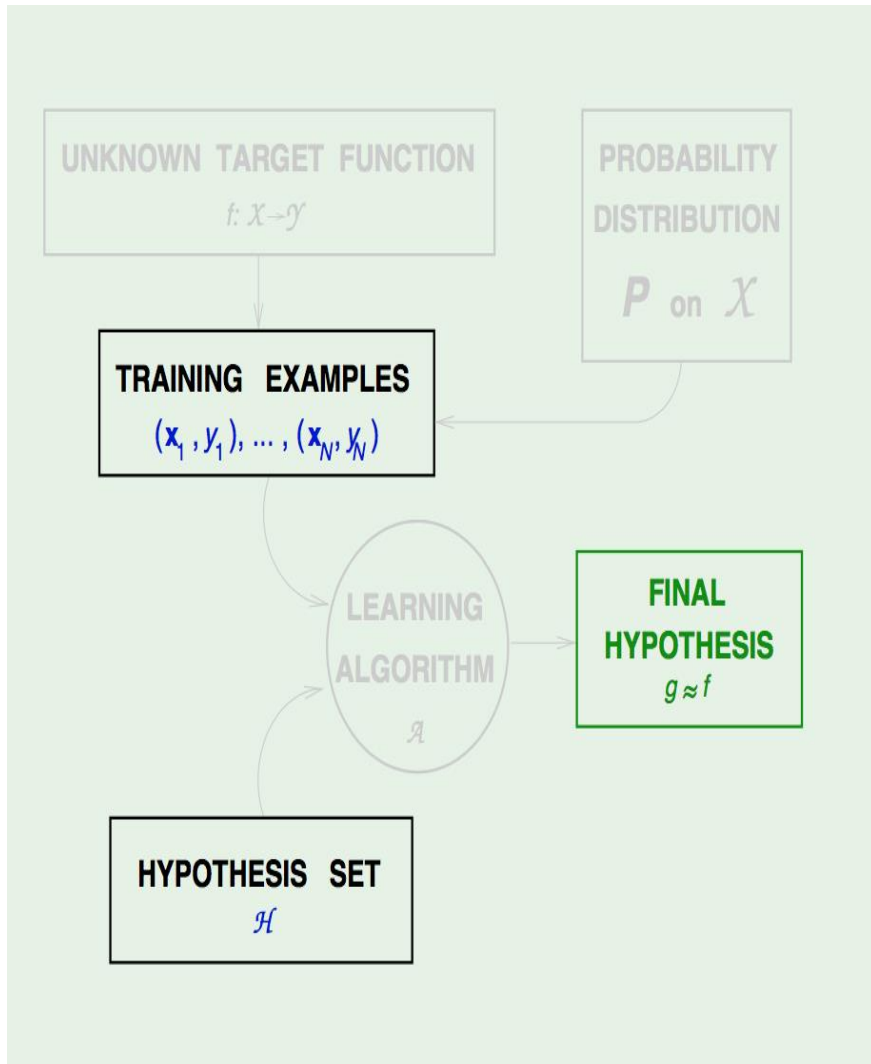
*With probability $\geq 1- \delta$*

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

$\mathcal{H}$ ~ model complexity

N: sample size

$1 - \delta$: confidence requirement

# 3. APPROXIMATION-GENERALIZATION [2]
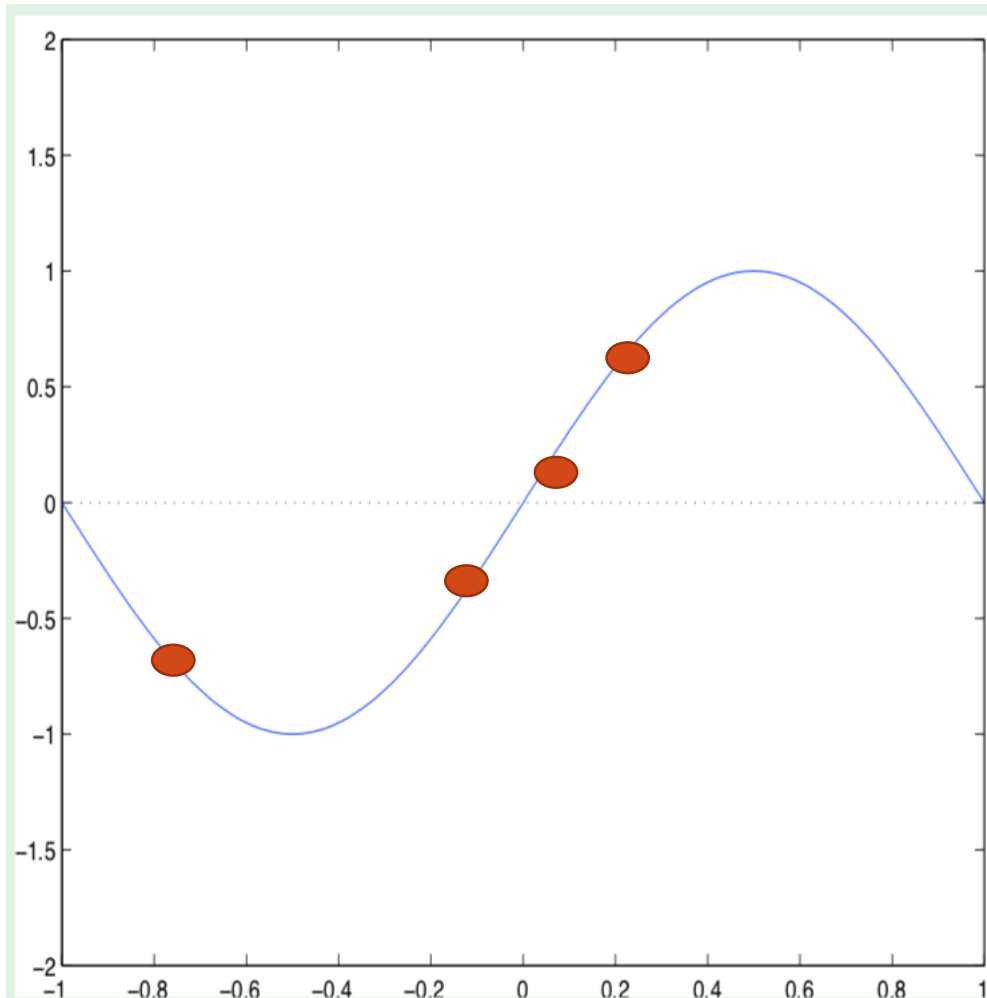


$$\mathbb{P}\big[|E_{\text{out}} - E_{\text{in}}| > \epsilon\big] \leq \underbrace{4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\epsilon^2 N}}_{\delta}$$

VC Dimension (1960 − 1990)
**"fundamental theory of learning"**
Vladimir Vapnik - Alexey Chervonenkis

$$\Omega(N, \mathcal{H}, \delta) = \sqrt{\frac{8}{N} \ln\left(\frac{4m_{\mathcal{H}}(2N)}{\delta}\right)}$$

Generalization bound

# 4. BIAS — VARIANCE TRADEOFF

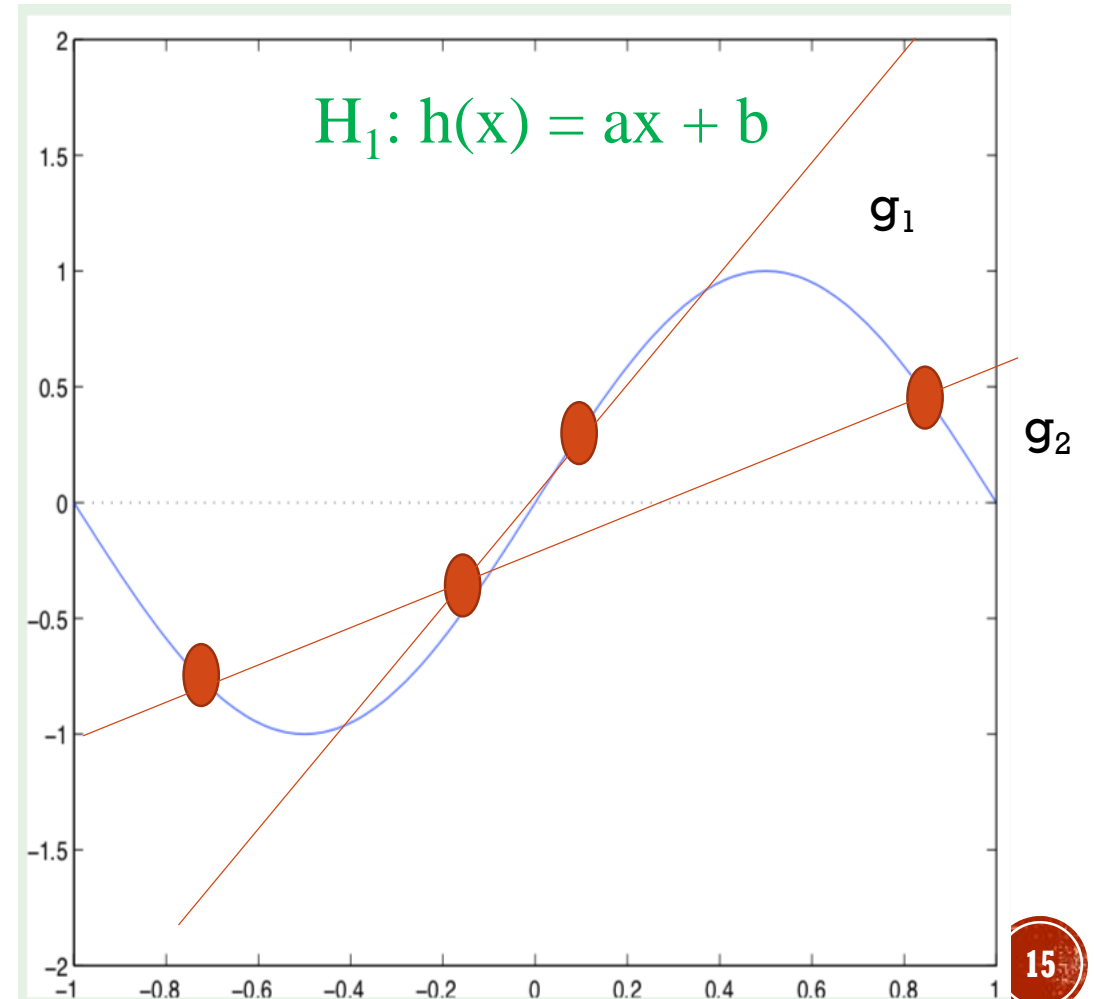

$$y = f(x) = \sin(\pi x).$$

$H_0: h(x) = b$    vs        $H_1: h(x) = ax + b$

## Which is better?
Approximation & generalization

# 4. BIAS — VARIANCE TRADEOFF



$H_0 \colon h(x) = b$

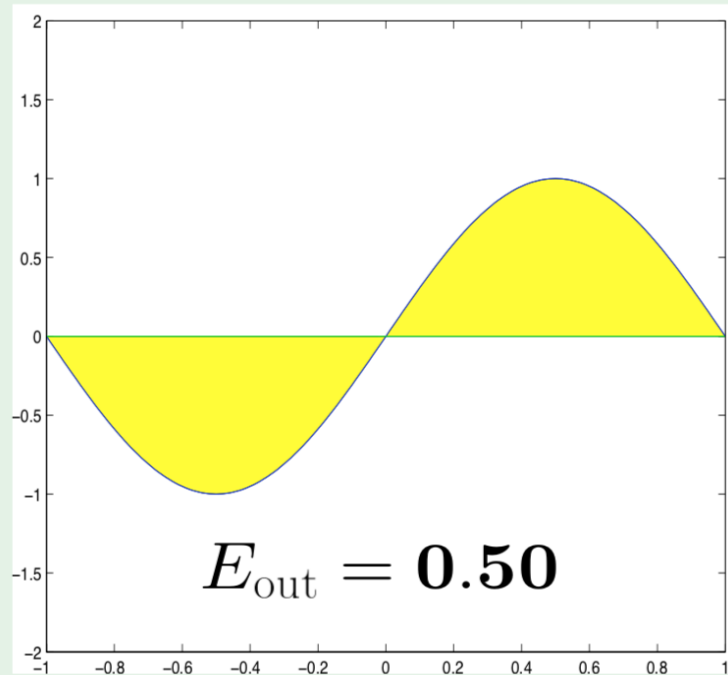$H_1 \colon h(x) = ax + b$

$g_1$
$g_2$

$g_1$
$g_2$

# 4. BIAS – VARIANCE TRADEOFF

"Approximation" - bias



$\mathcal{H}_0$        $\mathcal{H}_1$

$E_{\text{out}} = \mathbf{0.50}$        $E_{\text{out}} = \mathbf{0.20}$

16

# 4. BIAS – VARIANCE TRADEOFF

"Generalization" - Variance

# 4. BIAS – VARIANCE TRADEOFF [2][4]

Bias – Variance – who win ?



$\mathcal{H}_0$

$\bar{g}(x)$

$\sin(\pi x)$

bias = **0.50**   var = **0.25**

$\mathcal{H}_1$

$\bar{g}(x)$

$\sin(\pi x)$

bias = **0.21**   var = **1.69**

18

$$E_{out}\left(g^{(D)}\right) = \mathbb{E}_x\left[\left(g^{(D)}(x) - f(x)\right)^2\right]$$
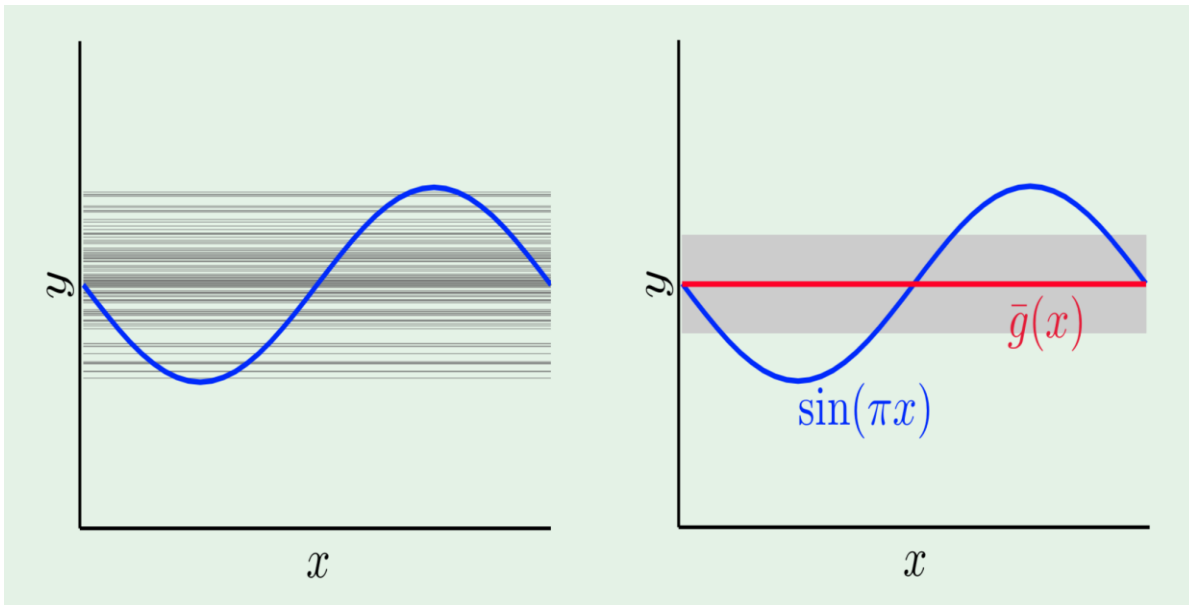
$$= \mathbb{E}_D\left[\left(g^{(D)}(x) - \bar{g}(x) + \bar{g}(x) - f(x)\right)^2\right]$$

$$= \cdots = \mathbb{E}_D\left[\left(g^{(D)}(x) - \bar{g}(x)\right)^2\right] + \mathbb{E}_D\left[\left(\bar{g}(x) - f(x)\right)^2\right]$$

19

# 3. BIAS - VARIANCE DECOMPOSITION [2]

$$E_{out}\left(g^{(D)}\right) = \mathbb{E}_D\left[\left(g^{(D)}(x) - \bar{g}(x)\right)^2\right] + \mathbb{E}_D\left[\left(\bar{g}(x) - f(x)\right)^2\right]$$

Bias                    Variance

Bias – variance decomposition $E_{out}$ to:

- How well H can approximate f

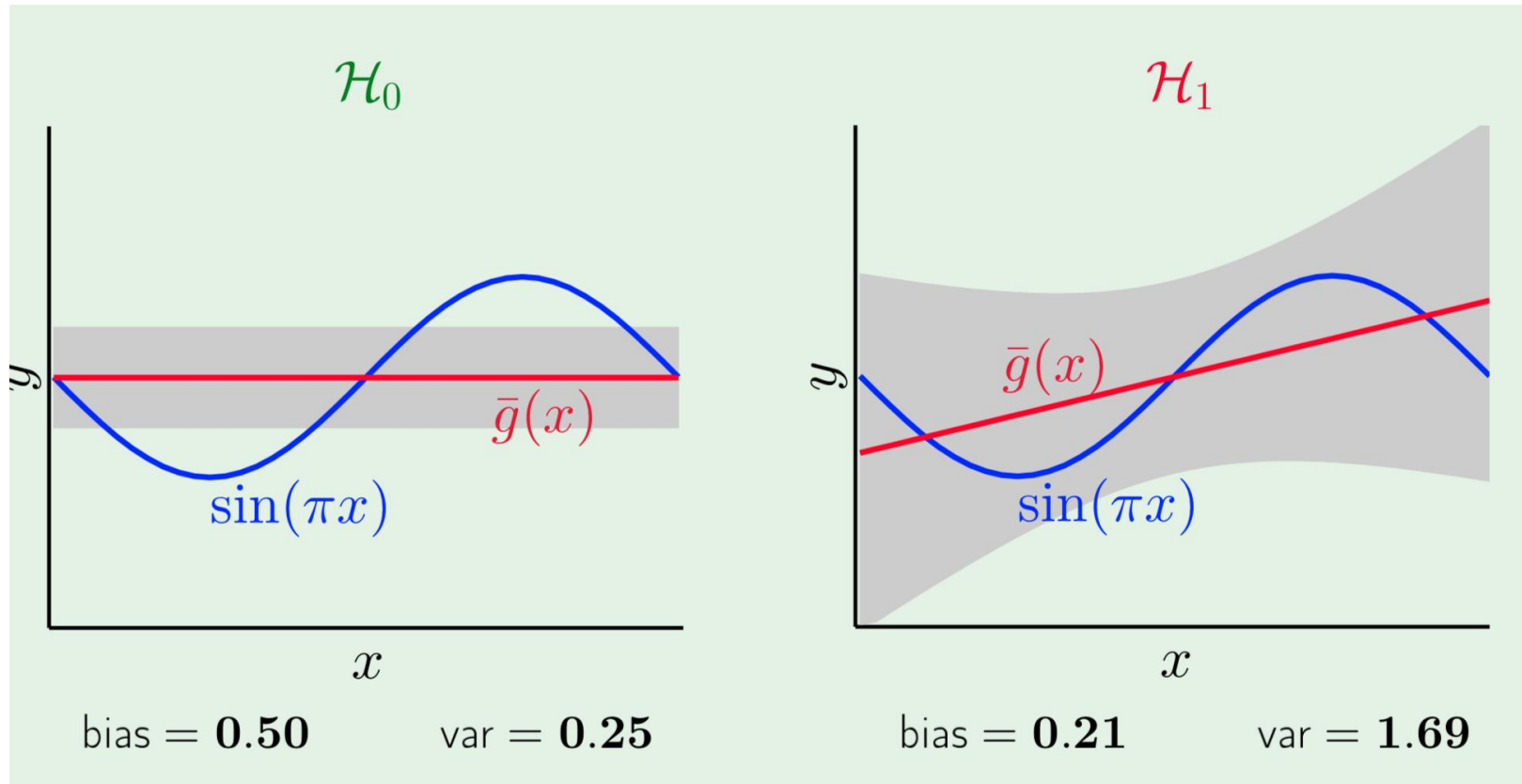- How well we can zoom in on a good h of H

Imagine **many** data sets $\mathcal{D}_1, \mathcal{D}_2, \cdots, \mathcal{D}_K$

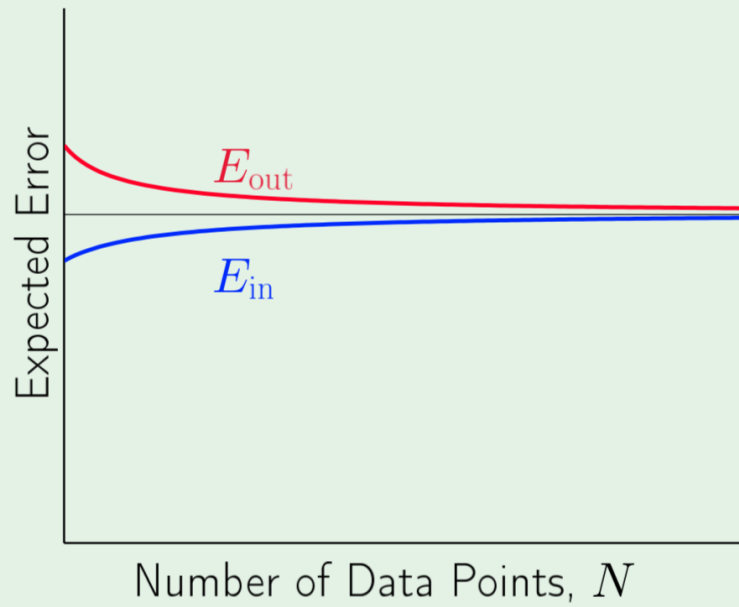$$\bar{g}(\mathbf{x}) \approx \frac{1}{K}\sum_{k=1}^{K}g^{(\mathcal{D}_k)}(\mathbf{x})$$

20

# 4. BIAS — VARIANCE TRADEOFF

WHO WON … ?        Congratulation $\mathcal{H}_0$

# 5. THE LEARNING CURVE

# BUT NOISE ....

# REFERENCES:

1. Pic: http://www.sthda.com/english/articles/40-regression-analysis/165-linear-regression-essentials-in-r/

2. Yaser S. Abu-Mostafa, Malik Magdon-Ismail, Hsuan-Tien Lin-Learning From Data. A short course-AMLBook (2012)

3. https://medium.com/@mp32445/understanding-bias-variance-tradeoff-ca59a22e2a83

4. Bishop - Pattern Recognition And Machine Learning - Springer  2006

# 3. GENERALIZATION [2]

$$\mathbb{P}\big[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\big] \leq 2Me^{-2\epsilon^2 N}$$
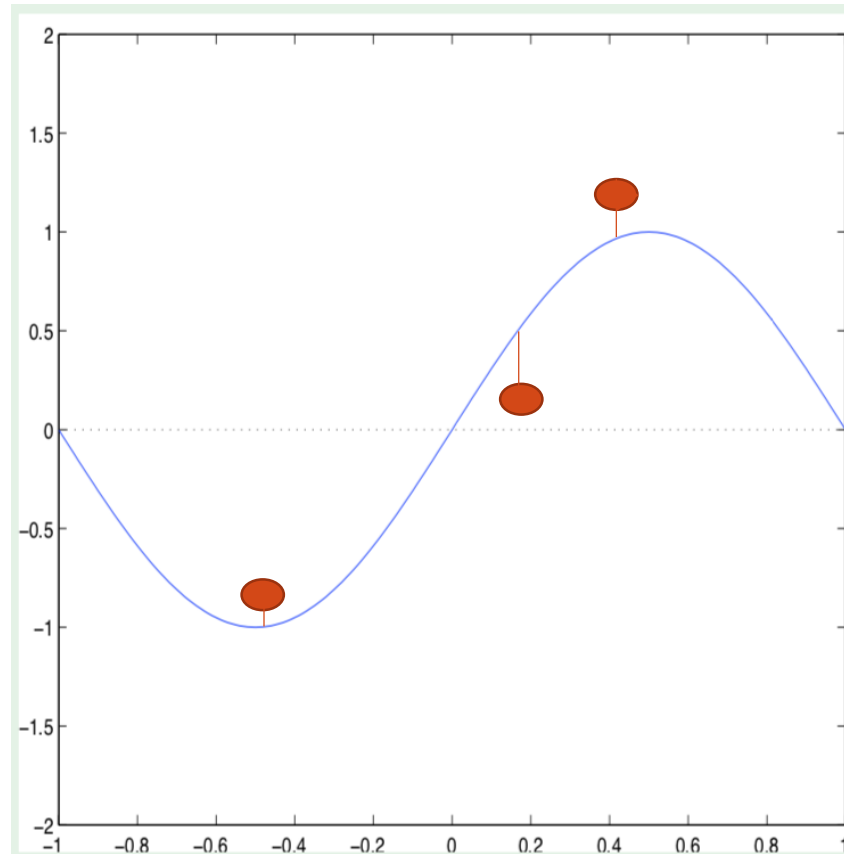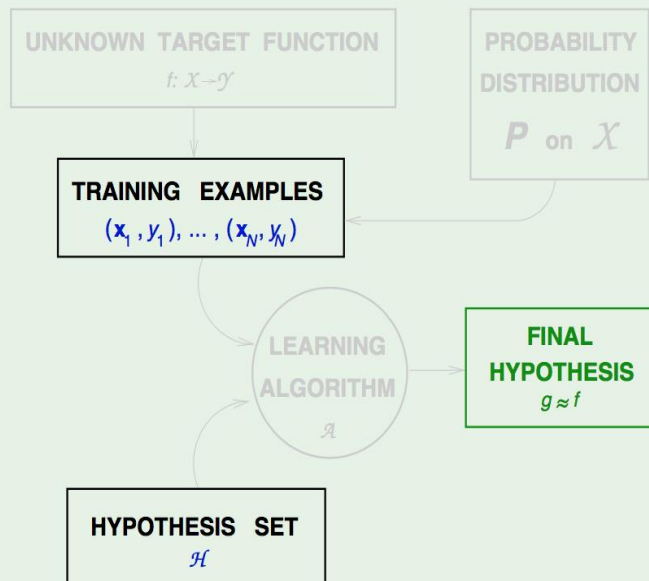
**VC Dimension (1960 – 1990)**
**"fundamental theory of learning"**
**Vladimir Vapnik - Alexey Chervonenkis**

$$\mathbb{P}\big[\,|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\,\big] \;\leq\; 4\; m_{\mathcal{H}}(2N)\; e^{-\frac{1}{8}\epsilon^2 N}$$

UNKNOWN TARGET FUNCTION
$f: \mathcal{X} \to \mathcal{Y}$

PROBABILITY
DISTRIBUTION
$P$ on $\mathcal{X}$

$$\mathbb{P}\big[|E_{\text{out}} - E_{\text{in}}| > \epsilon\big] \;\leq\; \underbrace{4 m_{\mathcal{H}}(2N) e^{-\frac{1}{8}\epsilon^2 N}}_{\delta}$$

TRAINING EXAMPLES
$(x_1, y_1), \ldots, (x_N, y_N)$

$$\delta = 4 m_{\mathcal{H}}(2N) e^{-\frac{1}{8}\epsilon^2 N} \implies \epsilon = \underbrace{\sqrt{\frac{8}{N} \ln \frac{4 m_{\mathcal{H}}(2N)}{\delta}}}_{\Omega}$$

LEARNING
ALGORITHM
$\mathcal{A}$

FINAL
HYPOTHESIS
$g \approx f$

HYPOTHESIS SET
$\mathcal{H}$

$$\text{With probability} \geq 1 - \delta, \qquad |E_{\text{out}} - E_{\text{in}}| \leq \Omega(N, \mathcal{H}, \delta)$$

# BIAS-VARIANCE VS VC.DIMENSION



VC analysis

bias-variance