

$$D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\} \subseteq X \times Y$$

$$X = \mathbb{R}^d$$

$$Y =$$

$\{0, 1\}$ ← binary classification

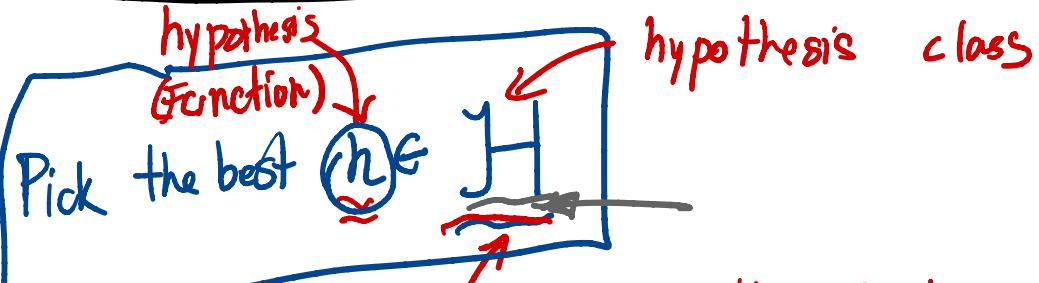
$\{0, 1, \dots, k-1\}$ ← multi-class classification

\mathbb{R} ← regression

Assume each $(\vec{x}_i, y_i) \in D$ is i.i.d. and is

drawn from distribution $P(\underline{(\vec{x}_i, y_i) \sim P})$.

Goal: Learn some function $\underline{h}: \underline{X} \rightarrow \underline{Y}$
such that $\underline{h}(\vec{x}) = y$ for any $(\vec{x}, y) \sim P$



The set of all possible functions we could learn

IDEAL: We wish to find the function h
that makes the Fewest mistakes

Loss functions

$$L(\underline{h}, D)$$

(percentage error)

- 0/1 Loss $L_{0/1}(h, D) = \frac{1}{|D|} \sum_{\forall (\vec{x}, y) \in D} d(h(\vec{x}), y),$

Binary classification:

+ Multi-class classification where

$$d(h(\vec{x}), y) = \begin{cases} 0 & \text{if } h(\vec{x}) = y \\ 1 & \text{if } h(\vec{x}) \neq y \end{cases}$$

- Square Loss $L_{sq}(h, D) = \frac{1}{|D|} \sum_{\forall (\vec{x}, y) \in D} (h(\vec{x}) - y)^2$

Regression:

Absolute loss

$$L_{ABS}(h, D) = \frac{1}{|D|} \sum_{\forall (\vec{x}, y) \in D} |h(\vec{x}) - y|$$

$$(|h(\vec{x}) - y|)$$

An ML algorithm is to choose $h^* \in \mathcal{H}$ that minimizes
the loss function

$$h^* = \arg \min_{h \in \mathcal{H}} L(h, D)$$

(local loss)

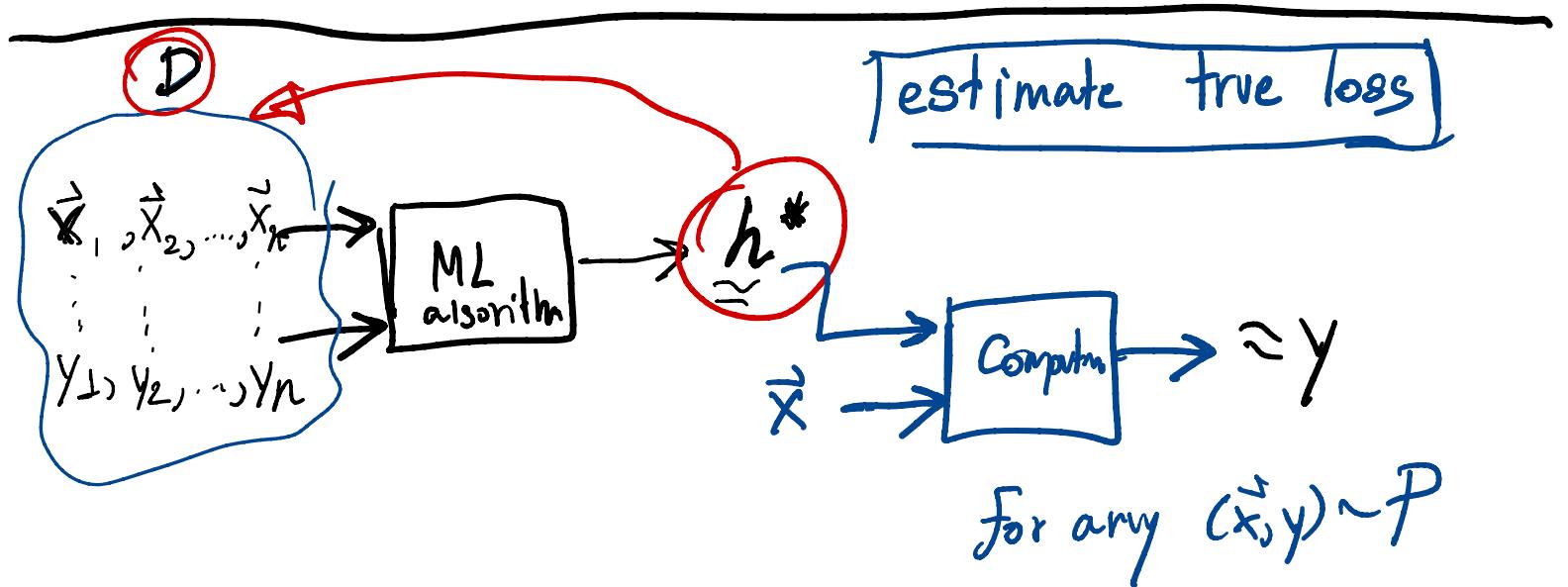
Exercise: Suppose we have the following h

$$h(\vec{x}) = \begin{cases} y_i & \text{if } \vec{x} = \vec{x}_i \text{ s.t. } \exists (\vec{x}_i, y_i) \in D \\ 0 & \text{otherwise} \end{cases}$$

Overfitting

100% D

Overfitting: We have the function h that works well on our data set but cannot be generalized to general data



True loss / Generalization loss: The true goal of learning

is to minimize

$$\mathcal{E}_g = \mathbb{E} \left[L(h, \{\vec{x}, y\}) \right] \quad (\vec{x}, y) \sim P$$

True goal: To learn h^* s.t.

$$h^* = \arg \min_{h \in \mathcal{H}}$$

$$\mathbb{E} \left[L(h, f(\vec{x}, y)) \right] \quad (\vec{x}, y) \sim P$$

Unknown

(Global loss)

Local loss:

$$L(h, D)$$

ML alg.
tries to minimize

Global loss:

$$\varepsilon = E[L(h, \{(\vec{x}, y)\})]$$

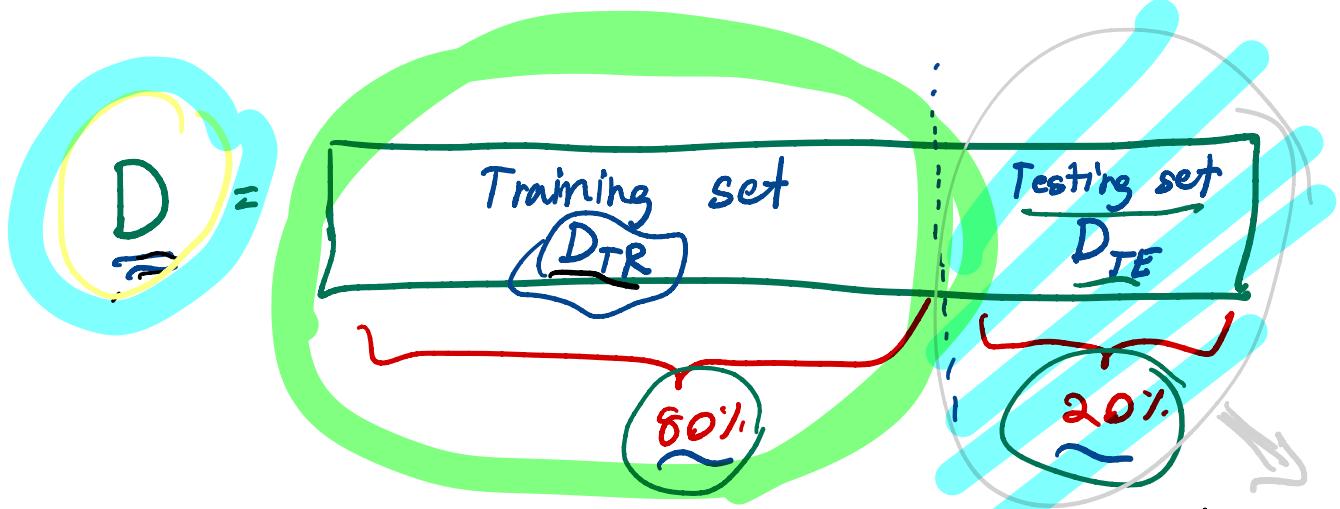


$$h(\vec{x}) \approx y$$

$$\forall (\vec{x}, y) \sim P$$

Our expectation
towards the
entire distribution

D represents the whole distribution P .



- ML algs. pick h^* s.t.

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} L(h, D_{TR})$$

- Evaluation of testing loss

$$\tilde{\epsilon}_{TE} = \tilde{L}(h^*, D_{TE})$$

find loss that
estimates the
global loss!

$$\tilde{\epsilon}_{TE} \approx \epsilon$$

$$D = D_{TB} \cup D_{TE}$$

$\frac{0.1}{10} \times 100\% = 1\%$

Claim: As $|D_{TE}| \rightarrow +\infty$, $\underline{\varepsilon}_{TE} \rightarrow \underline{\varepsilon}$.

Proof idea: $\lim_{|D_{TE}| \rightarrow +\infty} L(h, D_{TE}) = E[L(h, \{(x_i, y_i)\})]$

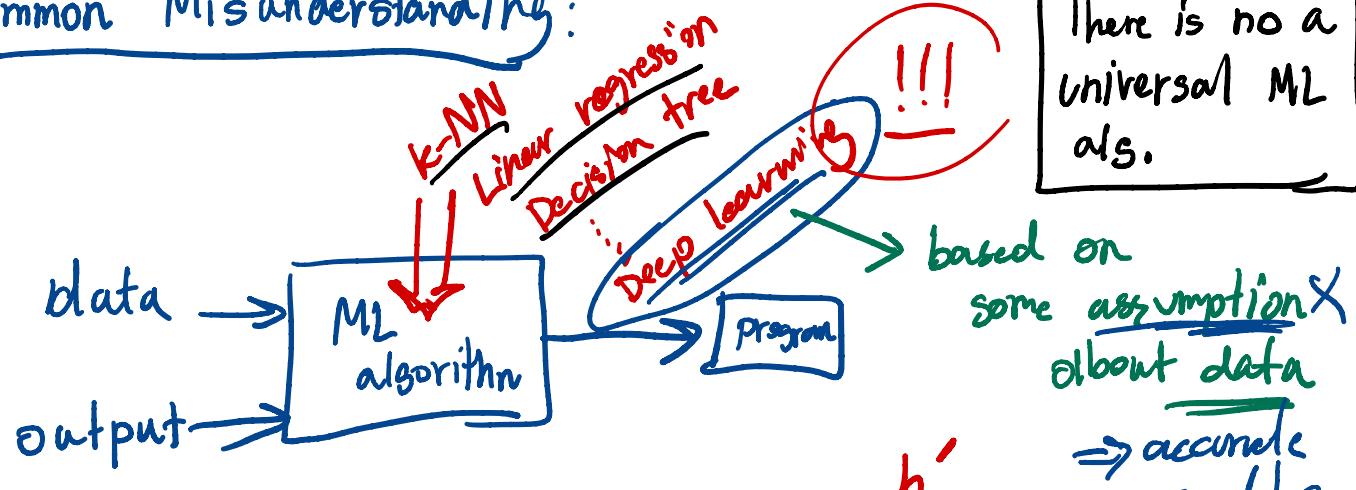
Therefore, $\hat{\epsilon}_{TE}$ is an unbiased estimator of ϵ .

The weak law of large number.

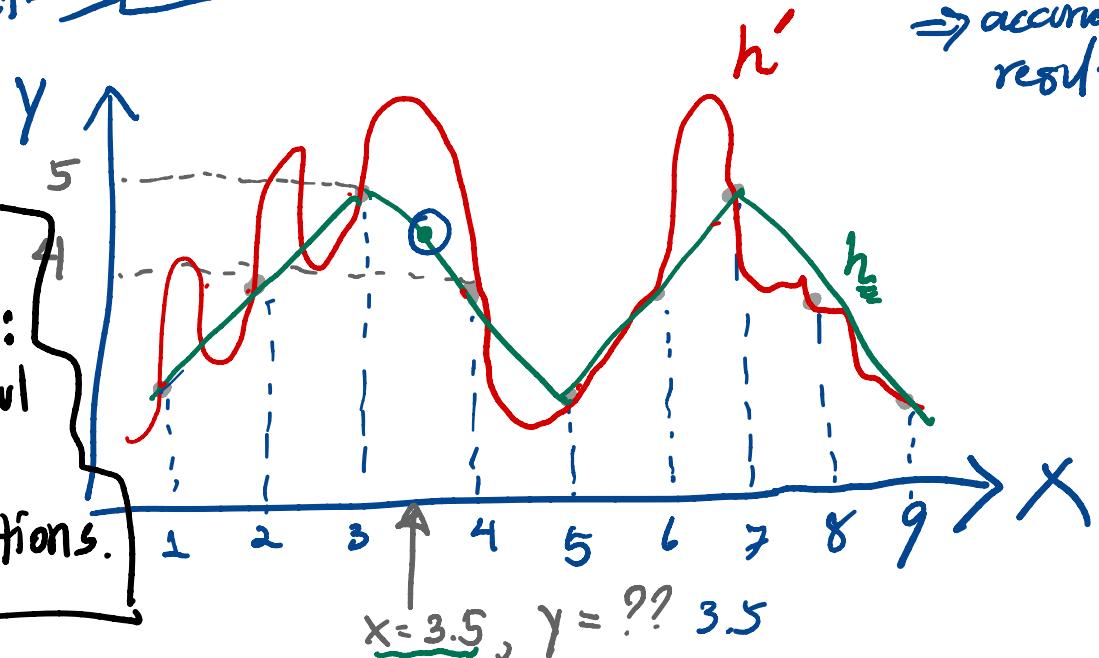
Summary: Supervised learning in practice.

- ① Define the setup
- ② Collect data
- ③ Training the program via training data
(choosing appropriate ML alg.)
- ④ Test the program via testing data.
(measure program's accuracy.)
- ⑤ The testing loss/error will indicate the actual accuracy of the program.

Common Misunderstanding:



No Free Lunch Theorem:
Every successful ML alg. must make assumptions.



$$L(h_3, D) = 0^2 + 0^2 + \frac{36}{(2-8)^2} + \frac{36}{(12-6)^2} + \frac{(8-4)}{16}$$

$$x=4 \rightarrow y=\underline{12} \quad | \quad h_3^{(4)} = \underline{8} \quad \frac{36+36}{+16}$$

$$x=8 \rightarrow y=12 \quad | \quad h_3^{(8)} = \underline{6} = \underline{\underline{88}}$$

$$x=12 \rightarrow y=\underline{8} \quad | \quad h_3^{(12)} = \underline{4}$$