

OPTIMAL LECTURE 1, 18/1/2021

ML converts experience in expertise/knowledge.

DATASET

Supervised learning

$$S = \{ (x_i, y_i) \mid i=1, \dots, n \}$$

$x \rightarrow ?$
 ↗ classification
 ↘ regression

unsupervised learning $S = \{ z_i, i=1, \dots, n \}$

clustering

reinforcement learning

$$\boxed{h: x \rightarrow y}$$

$$l(h, (x, y)) = (h(x) - y)^2$$

repr.

$$l(h, (x, y)) = \mathbb{1}(h(x) \neq y)$$

class.

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

HOW TO EVALUATE A PREDICTOR

tomorrow $(x, y) \sim D$

$$\underset{h \in H}{\text{minimize}} \quad \mathbb{E}_{(x, y) \sim D} [l(h, (x, y))]$$

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [l(h, (x, y))] \begin{cases} \rightarrow L_D(h) & \text{expected loss} \\ \rightarrow R(h) & \text{expected risk} \end{cases}$$

FIND A GOOD PREDICTOR \rightarrow SOLVE AN OPTIMIZATION PROBLEM

$$\underset{h \in \mathcal{H}}{\text{minimize}} \quad R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [l(h, (x, y))]$$

\mathcal{D} is unknown
ML is "distribution-free"

ML assumes S is drawn from \mathcal{D} (independently)

$$\frac{1}{n} \sum_{i=1}^n l(h, (x_i, y_i)) \approx \mathbb{E}_{(x,y) \sim \mathcal{D}} [l(h, (x, y))]$$

EMPIRICAL LOSS/RISK

$$L_S(h) / R_S(h)$$

$$(R_n(h))$$

CLOSER
AS $n \rightarrow \infty$

THE PROBLEM ML CONSIDERS IS

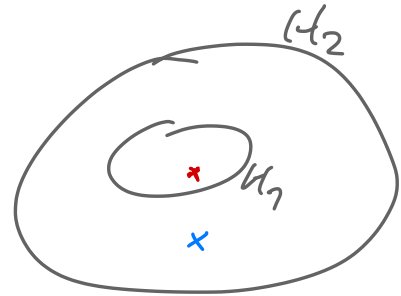
$$\underset{h \in H}{\text{minimize}} \quad R_S(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i))$$

choose between H_1, H_2

$$H_1 \subset H_2$$

$$h_1^* \in \underset{h \in H_1}{\text{argmin}} R_S(h)$$

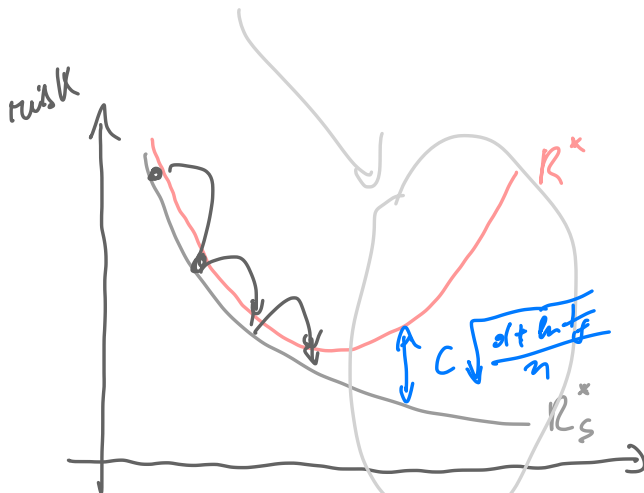
$$h_2^* \in \underset{h \in H_2}{\text{argmin}} R_S(h)$$



$$\underline{R_S(h_2^*) \leq R_S(h_1^*)}$$

OVERFITTING PROBLEM : IT MAY BE

$$R(h_2^*) > R_S(h_1^*)$$



Complexity of $H = d$

statistics

vs

\mathcal{H}

asymptotic
results

finite sample size

$$R_S(h) \simeq R(h)$$

$$|S|=n \rightarrow \infty$$

$\forall h \in \mathcal{H}$

$$|R_S(h) - R(h)| \leq C \sqrt{\frac{d + \ln \frac{1}{\delta}}{n}}$$

CONSTANT \leftarrow with prob. $1 - \delta$

$$d = \text{VC-dim}(\mathcal{H})$$

VC-dim

quantifies the complexity
of the class of predictors

for uni-dimensional
polynomials

$d = \#$ of coefficients
of the polynomial

Vapnik
Chervonenkis
dimension

$$R(h) \simeq R_S(h) + C \sqrt{\frac{d + \ln \frac{1}{\delta}}{n}}$$

if small it is ok to minimize
 $R_S(h)$ to minimize $R(h)$

REGULARIZATION

TERM

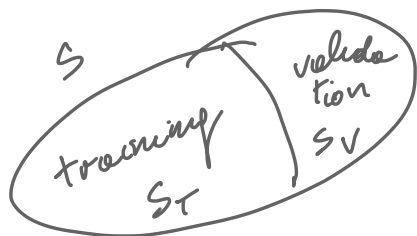
$$\text{minimize } R_S(h) + \lambda \|h\|^2$$

WE COULD

$$\text{minimize } R_S(h) + C \sqrt{\frac{d + \ln \frac{1}{\delta}}{n}}$$

BUT IT IS A LOOSE BOUND

IN PRACTICE:



$$h_1^* \in \arg\min_{h \in H_1} R_{S_T}(h)$$

$$h_2^* \in \arg\min_{h \in H_2} R_{S_T}(h)$$

$$(R_{S_T}(h_1^*) \geq R_{S_T}(h_2^*))$$

$$R_{S_V}(h_1^*) < R_{S_V}(h_2^*) \Rightarrow \text{SELECT } h_1^*$$

$$R_{S_V}(h_1^*) > R_{S_V}(h_2^*) \Rightarrow \text{SELECT } h_2^*$$

$$R(h_1^*) \stackrel{?}{\geq} R(h_2^*)$$

$$R(h_1^*) \simeq R_{S_V}(h_1^*)$$

$$R(h_2^*) \simeq R_{S_V}(h_2^*)$$

$$\left| R(h) - R_{S_T}(h) \right| \leq C \sqrt{\frac{VC\text{-dim}(H) + \ln \frac{1}{\delta}}{|S_T|}} \quad VC\text{-dim}(H) \text{ BIG!}$$

$\forall h \in H$

I only care of h_1^* & h_2^*

$$\left| R(h) - R_{S_V}(h) \right| \leq C \sqrt{\frac{VC\text{-dim}(\{h_1^*, h_2^*\}) + \ln \frac{1}{\delta}}{|S_V|}}$$

$h \in \{h_1^*, h_2^*\}$

$$VC\text{-dim}(\{h_1^*, h_2^*\}) \ll VC\text{-dim}(H)$$

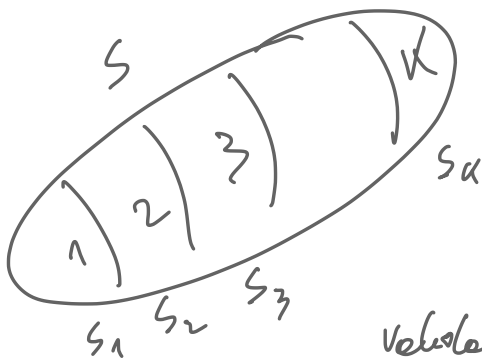
$$\hookrightarrow \leq 1$$

IT CAN BE

$R_{S_T}(h_1^*)$ is not close to $R(h_1^*)$

but $R_{S_V}(h_1^*)$ is close to $R(h_1^*)$

k-fold cross validation



$$h_1 \in \operatorname{argmin}_{S_2 \cup \dots \cup S_k} R(h)$$

$$h_2 \in \operatorname{argmin}_{S_1 \cup S_3 \cup \dots \cup S_k} R(h)$$

$$h_k \in \operatorname{argmin}_{S_1 \cup S_2 \cup \dots \cup S_{k-1}} R(h)$$

$$\text{validation} = \frac{1}{k} \sum_{i=1}^k R_{S_i}(h_i)$$

K cross fold validation \rightarrow solving K optimization problem

HYPER PARAMETERS

= PARAMETERS OF THE ALGORITHM
USED TO LEARN THE MODEL
 \neq PARAMETERS OF THE MODEL

- learning rate
- parameters of the grid search
- batch size

EXPERIENCE : $3 \div 10$ VALUES PER HYPER PARAMETER

LEARNING PROBLEM

models $\times K$ $\times (3 \div 10)^{\text{\# hyperparameters}}$
(K -fold)

