

Learning passive

Best strategy of rats: Rats eat food, get sick, and afterwards avoid the area for a while. This comes up with the target value.

Sitting  
 - supervised: We have from labeled training data with  $x$  to get values. should try to derive a rule that comes up with the target value.  
 - passive learner: The learner has no influence on the environment and its internal state of space.

Training process  
 - batch learning: target are labeled training data. Set  $S$  and we have to learn from it and not to predict on that afterwards.  
 How to formalize this?

- Domain  $X$ : The set of all possible examples  $x \in \mathbb{R}^n$
- Label set  $Y$ : The set of possible target values, often  $Y = \{0, 1\}$
- example complexity: We have a distribution  $\mathcal{D}$  over  $X \times Y$ . There is some loss  $\ell$ .  $\mathcal{D}$  has a measure of entropy.

$$L_{\mathcal{D}}(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}}[\ell(h(x), y)]$$

Learning Task

Input: A training sequence  $S = ((x_1, y_1), \dots, (x_m, y_m))$  drawn independently and identically distributed from  $\mathcal{D}$

Output: a prediction rule  $A(S)$  that can, given  $x \in X$  produce a value  $A(S)(x) \in Y$  s.t.  $L_{\mathcal{D}}(A(S))$  is minimized

First learning rule: Empirical Risk Minimization (ERM)

Given a training sequence  $S$ , return a prediction rule that minimizes the empirical error (on the training sequence)

$$L_S(A(S)) = \frac{|\{i \in [m] : A(S)(x_i) \neq y_i\}|}{m}$$

How to fix the problem of memorizing / overfitting?  
 We introduce a restricted search space  $\mathcal{H}$ , called hypothesis class.

Then ERM becomes the following:

Given:  $S$  and some finite representation of  $\mathcal{H}$

return:  $A(S) = \text{ERM}_{\mathcal{H}}(S) \in \arg \min_{h \in \mathcal{H}} L_S(h)$

notation: more restricted  $\mathcal{H} \Rightarrow$  a better signal to noise ratio  $\Rightarrow$  stronger inductive bias

PAC & APAC  
 Which hypothesis classes are learnable?  
 we can come up with an algorithm that for any distribution  $\mathcal{D}$  and training sequence  $S$  minimizes  $L_{\mathcal{D}}(A(S))$  with some guarantee.

DEF: A hypothesis class  $\mathcal{H}$  (for some  $X, Y$ ) is APAC learnable if there exists:  
 - a function  $m_{\mathcal{H}}(\epsilon, \delta) \rightarrow \mathbb{N}$   
 - a learning algorithm  $A$  with:

$\forall \epsilon, \delta \in (0, 1), \forall$  distributions  $\mathcal{D}$  over  $X \times Y$   
 ~~$\forall$  labeling functions  $f: X \rightarrow Y$~~   
~~with test error  $\exists h \in \mathcal{H}$  s.t.  $L_{\mathcal{D}}(h) = 0$~~   
 thresholds: when using  $A$  on  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  i.i.d. examples from drawn from  $\mathcal{D}$  labeled by  $f$ , the algorithm returns  $A(S) \in \mathcal{H}$  (with prob.  $1 - \delta$ )  $L_{\mathcal{D}}(A(S)) \leq \epsilon + \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$

Corr: Every realizable finite hypothesis class  $\mathcal{H}$  is APAC learnable using ERM with sample complexity  $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2 \log(|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil$

No Free Lunch

Let  $\mathcal{H} \subseteq \{f: X \rightarrow \{-1, 1\}\}$   
 Let  $A$  be any algorithm for binary classification s.t. to 0-1 loss over  $X$ . Let  $m \leq \frac{|X|}{2}$  (fix) then there exists a distribution  $\mathcal{D}$  over  $X \times Y$  s.t.  
 1)  $\exists f: X \rightarrow Y$  with  $L_{\mathcal{D}}(f) = 0$   
 2) With prob. at least  $\frac{1}{7}$  over  $S \sim \mathcal{D}^m$  we have  $L_{\mathcal{D}}(A(S)) \geq \frac{1}{8}$

DEF: Shattering

A hyp. class  $\mathcal{H}$  shatters some finite set  $C \subseteq X$  if  $\mathcal{H}|_C$  is the set of all functions from  $C$  to  $\{0, 1\}$

DEF: VC-Dimension

The VC-dimension of a hyp. class  $\mathcal{H}$  is the cardinality of the largest set  $C \subseteq X$  that is shattered by  $\mathcal{H}$