

- What is the inductive bias in machine learning? Reducing the search space for \mathcal{H}
- What does it mean that machine learning is distribution-free? It works well over X
- Definition of PAC-learning If PAC if under some and 0-1 loss $\exists A, \mathcal{H}_{\text{size}}(A) \leq M$, $\forall \epsilon > 0, \forall \delta > 0$ $\exists n$ such that $\Pr_{x \sim D}[\text{err}(h_n) \geq \epsilon] \leq \delta$
- What is the ERM? $\arg \min_{h \in \mathcal{H}} L_\Delta(h)$
- What is the realizability assumption? \mathcal{H}
- What is the sample complexity? Minimum number of samples of x w.r.t. \mathcal{H} for which $A(\mathcal{H})$ can efficiently learn $\log(M/\delta)$
- Are finite hypothesis classes learnable? What's their sample complexity? Yes, $\frac{\log(M/\delta)}{\epsilon^2}$
- Differences between Agnostic PAC learnability and PAC learnability In PAC for $\epsilon = 0$ $L_\Delta(h) = 0$ because $\Pr_{x \sim D}[h(x) \neq f(x)] = 0$
- Definition of Agnostic PAC learnability
- What is in general the expected loss of a predictor? $L_\Delta(h) = \mathbb{E} \left[\sum_i \ell(h(x_i), y_i) \right]$
- What is uniform convergence? UC: $\exists M_{\text{UC}}(\epsilon, \delta) \rightarrow N$ such that $\Pr_{x \sim D}[\text{err}(h_N) \geq \epsilon] \leq \delta$ when $N \geq M_{\text{UC}}(\epsilon, \delta)$
- If a hypothesis class has the uniform convergence property, what can you tell me about its learnability? UC $\Rightarrow \mathcal{H}$ is representable $\Rightarrow H \text{ PAC} \Rightarrow H \text{ PAC}$
- What is the sample complexity of a finite hypothesis class? $\frac{\log(M/\delta)}{\epsilon^2}$
- What is the discretization trick?
- Are infinite classes PAC-learnable? And Agnostic PAC-learnable?
- What are the consequences of the No Free Lunch theorem? $\mathbb{E}[\text{err}] = \infty$ when $H_{\text{size}}(h) \rightarrow \infty$ for any ϵ -loss
- What is the bias-complexity trade-off? $\text{Err}_{\text{app}}(H) \leq \text{Err}_{\text{est}}(H)$ $\text{Err}_{\text{est}}(H) \leq \text{Err}_{\text{app}}(H)$
- What is the approximation error? $L_\Delta(A(\mathcal{H})) = L_\Delta(A(\mathcal{H})) - \min_{h \in \mathcal{H}} L_\Delta(h) + \min_{h \in \mathcal{H}} L_\Delta(h)$
- What is the estimation error? $L_\Delta(A(\mathcal{H})) = L_\Delta(A(\mathcal{H})) - \min_{h \in \mathcal{H}} L_\Delta(h) + \min_{h \in \mathcal{H}} L_\Delta(h)$
- Advantages/disadvantages of selecting a more complex hypothesis class Err_{app} decreases but Err_{est} increases.
- What is the Bayes error? $L_\Delta(A(\mathcal{H})) = L_\Delta(A(\mathcal{H})) - \min_{h \in \mathcal{H}} L_\Delta(h) + \min_{h \in \mathcal{H}} L_\Delta(h) - \min_{h \in \mathcal{H}} L_\Delta(h)$
- Definition of shattering H shatters $A \subseteq X$ if $\forall h \in H, \Pr_{x \sim D}[\text{err}(h) = 0] = 1$
- VC dimension $V(H) = \max_{\mathcal{X} \subseteq X} |\{x_1, x_2, \dots, x_n\}|$ such that $\Pr_{x \sim D}[\text{err}(h) = 0] = 1$
- VC dimension and size of the class \mathcal{H} $\leq \text{VC}(H)$ then $V(H) \leq \text{VC}(H)$
- What if VC-dimension is infinite? $V(H) = \infty \Rightarrow H \text{ not PAC}$
- ERM is a successful PAC-learner for a class H , what consequences? $\text{Err}_{\text{app}}(A(\mathcal{H})) = 0$
- ERM is a successful agnostic PAC-learner for a class H , what consequences?
- What is the sample complexity of a PAC learnable class? $\approx \frac{V(H) + \log(1/\delta)}{\epsilon^2}$ by QF TSL
- What is the sample complexity of an agnostic PAC learnable class? $\approx \frac{1}{\epsilon^2}$
- How many samples are needed for a class to have the Uniform Convergence property? $M \geq M_{\text{UC}}(\epsilon, \delta) \approx \frac{d \cdot \ln(1/\delta)}{\epsilon^2}$
- Good points and negative points of ERM A finds the best predictor, but for H too complex can be computationally hard
- What is the problem with ERM? A can be NP-hard
- Define a linear classifier $h(x) = \sum_{i=1}^n w_i \cdot x_i + b$
- Are linear classifiers (PAC/Agnostic PAC) learnable? Why? How? Yes. In the separable case $V(H) = \infty$ so we can find the set of parameters by LP or perceptron learning rule
- Define linear regression
- Computational complexity of solving a linear classification problem with the 0-1 loss function $\text{Pr}_{x \sim D}[\text{err}(h) \geq 1]$ and $\log(1/\delta)$ in $\mathcal{O}(n^2)$
- Definition of weak learnability $\text{WL} \subseteq \text{PAC}$ $\forall A, \forall \text{f}, \forall \epsilon > 0, \exists n \in \mathbb{N}, \Pr_{x \sim D}[\text{err}(h_n) \geq \epsilon] \leq \delta$ when $(1-\epsilon) \cdot m \geq M_{\text{UC}}(\epsilon, \delta)$
- Are there classes that are weak learnable but not PAC learnable? $\text{PAC} \not\subseteq \text{WL}$ $\text{L}_\Delta(A(\mathcal{H})) \leq \frac{1}{2} - \epsilon$ up to ϵ w.r.t. def of WL
- Why have we discussed weak learnability? By boosting algorithms we can turn up weak A for B to find a strong B
- Describe the AdaBoost algorithm $\text{AdaBoost}(\mathcal{H}, T)$ for $i=1 \dots T$
 $h_i = \text{WL}(\mathcal{D}_i)$
 $\mathcal{D}_i = \mathcal{D}_{i-1} \setminus \{(x_i, y_i)\}$
 $w_i = \frac{1}{T} \cdot \ln \left(\frac{1 - \text{err}(h_i)}{\text{err}(h_i)} \right)$
 $\mathcal{D}_{i+1} = \mathcal{D}_i \cup \{(\mathbf{x}_i, \mathbf{y}_i)\}$
 $\text{Not correctly classified samples}$
- Guarantees for the AdaBoost algorithm
- What is the effect of the number of iterations in the AdaBoost algorithm?

If WL returns a classifier h then $L_\Delta(h_{\text{Ada}}) = e^{-2T} \cdot \mathbb{W}^T (1 - f_T)$

43. When is a set convex?
 44. What is a convex function?
 45. How can you check that a function is convex?
 46. Why do we care of convexity in ML?
 47. When is a ML problem convex? $\nabla \mathcal{L}(w(x, y))$
 48. Linear regression
 49. What are surrogate loss functions? convex function from the approx. more convex on \mathcal{H}
 50. Logistic regression
 51. What is the optimization error?
 52. Definition of non-uniform learnability
 53. When is a class non-uniform learnable?
 54. What is Structural Risk minimization?
 55. What is the minimum description length?
 56. What approaches do you know to select the model/the class of hypothesis class?
 57. How can we evaluate the quality of a trained model?
 58. What is the validation approach to model selection?
 59. What is the train-validation-test split?
 60. What is k-fold cross validation?
 61. What is the model selection curve?
 62. What can we do if we are not happy with the quality of the model we trained?

$$\mathcal{L}_D(\mathcal{A}(S)) = \mathcal{L}_D(\mathcal{A}(S)) + \mathcal{L}_{\text{err}}^{\text{sur}}(\mathcal{A}(S)) - \min_{h \in \mathcal{H}} \mathcal{L}_{\text{err}}^{\text{sur}} \min_{h \in \mathcal{H}} \mathcal{L}_{\text{err}}^{\text{sur}} = m/m \mathcal{L}_{\text{err}}^{\text{sur}} + m/m \mathcal{L}_{\text{err}}^{\text{sur}}$$

H NUC

52) H NUC $\Leftrightarrow \exists A: (x, y) \rightarrow H \exists m_{\text{uc}}: (0, 1) \times H \rightarrow \mathbb{N}, \forall \delta \in (0, 1), \forall h \in H$
 $S \text{ iid sample of } D \text{ of size } m \geq m_{\text{uc}}(\epsilon, \delta) \text{ samples the}$
 $\mathcal{L}_D(\mathcal{A}(S)) \leq \mathcal{L}_S(h) + \epsilon \text{ w.p. } 1 - \delta$

53) H is NUC $\Leftrightarrow H = \bigcup_m \mathcal{H}_m / \text{VCdim}(\mathcal{H}_m) < +\infty$

Structural Risk Minimization (SRM)

prior knowledge:

$\mathcal{H} = \bigcup_n \mathcal{H}_n$ where \mathcal{H}_n has uniform convergence with $m_{\mathcal{H}_n}^{\text{uc}}$
 $w: \mathbb{N} \rightarrow [0, 1]$ where $\sum_n w(n) \leq 1$

define: ϵ_n as in Equation (7.1); $n(h)$ as in Equation (7.4)

input: training set $S \sim \mathcal{D}^m$, confidence δ

output: $h \in \arg\min_{h \in \mathcal{H}} [L_S(h) + \epsilon_{n(h)}(m, w(n(h)) \cdot \delta)]$

54) 2 ways:

1) STRUCTURAL RISK MINIMIZATION

$$\begin{aligned} &\text{Let } H = \bigcup_m \mathcal{H}_m \text{ has UC} \\ &\text{Let } w: \mathbb{N} \rightarrow [0, 1] \mid \sum_n w_n = 1 \\ &\text{Let } \epsilon_m(\delta) = \min \left\{ \epsilon \mid \mathbb{E} \left[\sum_{h \in \mathcal{H}_m} \mathbb{I}_{\{h \text{ fails}\}} \right] \leq \delta \right\} \Rightarrow \min \left\{ \epsilon \mid \frac{\text{VCdim}(\mathcal{H}_m) \cdot \log(m/\delta)}{m} \leq \epsilon \right\} \\ &\text{Let } n(h) = \min \left\{ n(h) \mid h \in \mathcal{H}_m \right\} \\ &\text{output h} = \arg\min_{h \in \mathcal{H}} L_S(h) + \epsilon_m(m, w_m \cdot \delta) \Rightarrow h = \arg\min_{h \in \mathcal{H}} L_S(h) + \sqrt{\frac{\text{VCdim}(\mathcal{H}_m) \cdot \log(m/\delta)}{m}} \end{aligned}$$

2) VALIDATION SET To choose among $\mathcal{H}_1, \dots, \mathcal{H}_m$ we compute $E_{\mathcal{H}_i}, i=1 \dots n$

Then we must compare $L_D(h_1), \dots, L_D(h_m)$ to choose

Not having D we use the VALIDATION SET choosing $\arg\min_{\{h_i\} \subseteq \mathcal{H}} L_{\text{val}}(h_i)$

We formally check

$$|L_{\text{val}}(h_i) - L_D(h_i)| \leq \sqrt{\frac{\Omega_m(2/\delta)}{2m_{\text{val}}}}$$

$$\sqrt{\sum_{h \in \mathcal{H}_i} \mathbb{P}(h \text{ fails})}$$

S7) TS AND TESTS

- given $h = \arg \min_h L_{TS}(h)$

. check

$$L_D(h) \approx L_{tests}(h) = \frac{1}{m_{tests}} \sum_{(x_i, y_i) \in tests} I(h(x_i) \neq y_i)$$

. we shall have

$$|L_{tests}(h) - L_D(h)| \leq \epsilon_{test} = \sqrt{\frac{9m(2k)}{m_{test}}}$$

Q) When data is scarce how we apply validation without waste of data?

```

for i=1 - m      {u_i, - h_i}
  {
    for j=1 - k      {s_j, - s_k}
      {
        h_ij = ORTH_{u_i}(s_j)
        f_{ij} = Loss(h_ij)
        q_i = 1/k * sum_{j=1}^k f_{ij}
      }
    i* = argmin_{i=1 - m} q_i
    h = BRM_{h*}(S)
  }

```