

RECITATION 10: LEARNING THEORY AND ENSEMBLE METHODS

10-301/10-601 Introduction to Machine Learning (Summer 2024)
<http://www.cs.cmu.edu/~hchai2/courses/10601>

1 Learning Theory

1.1 PAC Learning

Some Important Definitions

Basic notation:

- Probability distribution (unknown): $X \sim p^*$
- **True function** (unknown): $c^* : X \rightarrow Y$
- **Hypothesis space** \mathcal{H} and **hypothesis** $h \in \mathcal{H} : X \rightarrow Y$
- Training dataset $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$

2. True Error (expected risk)

$$R(h) = P_{x \sim p^*}(c^*(x) \neq h(x))$$

3. Train Error (empirical risk)

$$\begin{aligned}\hat{R}(h) &= P_{x \sim \mathcal{D}}(c^*(x) \neq h(x)) \\ &= \frac{1}{N} \sum_{i=1}^N 1(c^*(x^{(i)}) \neq h(x^{(i)})) \\ &= \frac{1}{N} \sum_{i=1}^N 1(y^{(i)} \neq h(x^{(i)}))\end{aligned}$$

The **PAC criterion** is that we produce a high accuracy hypothesis with high probability. More formally,

$$P(\forall h \in \mathcal{H}, \text{_____} \leq \text{_____}) \geq \text{_____}$$

Sample Complexity is the minimum number of training examples N such that the PAC criterion is satisfied for a given ϵ and δ

Sample Complexity for 4 Cases: See Figure 1. Note that

- **Realizable** means $c^* \in \mathcal{H}$
- **Agnostic** means c^* may or may not be in \mathcal{H}

	Realizable	Agnostic
Finite $ \mathcal{H} $	Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.
Infinite $ \mathcal{H} $	Thm. 3 $N = O(\frac{1}{\epsilon} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 4 $N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.

12

Figure 1: Sample Complexity for 4 Cases

The **VC dimension** of a hypothesis space \mathcal{H} , denoted $\text{VC}(\mathcal{H})$ or $d_{\text{VC}}(\mathcal{H})$, is the maximum number of points such that there exists at least one arrangement of these points and a hypothesis $h \in \mathcal{H}$ that is consistent with any labelling of this arrangement of points.

To show that $\text{VC}(\mathcal{H}) = n$:

-
-

Questions

1. For the following examples, write whether or not there exists a dataset with the given properties that can be shattered by a linear classifier.
 - 2 points in 1D
 - 3 points in 1D
 - 3 points in 2D
 - 4 points in 2D

How many points can a linear boundary (with bias) classify exactly for d-Dimensions?

2. Consider a rectangle classifier (i.e. the classifier is uniquely defined 3 points $x_1, x_2, x_3 \in \mathbb{R}^2$ that specify 3 out of the four corners), where all points within the rectangle must equal 1 and all points outside must equal -1

(a) Which of the configurations of 4 points in figure 2 can a rectangle shatter?

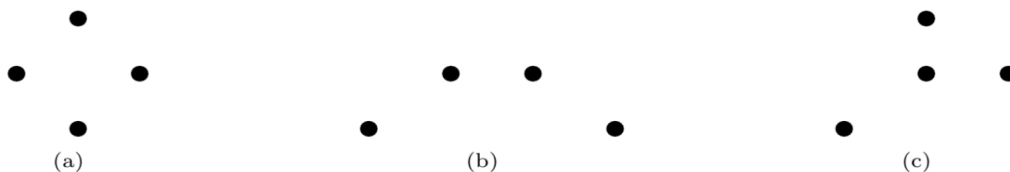


Figure 2

(b) What about the configurations of 5 points in figure 3?



Figure 3

3. Let x_1, x_2, \dots, x_n be n random variables that represent binary literals ($x \in \{0, 1\}^n$). Let the hypothesis class \mathcal{H}_n denote the conjunctions of no more than n literals in which each variable occurs at most once. Assume that $c^* \in \mathcal{H}_n$.

Example: For $n = 4$, $(x_1 \wedge x_2 \wedge x_4), (x_1 \wedge \neg x_3) \in \mathcal{H}_4$

Find the minimum number of examples required to learn $h \in \mathcal{H}_{10}$ which guarantees at least 99% accuracy with at least 98% confidence.

2 Ensemble Methods

The idea of ensemble methods is to build a model for prediction by combining the strengths of a group of simpler models. We'll cover two examples of ensemble methods: random forests and AdaBoost.

2.1 Random Forests

1. What are some downsides of decision trees, and how can we explain this in the context of the bias-variance tradeoff?

Random Forests = Sample Bagging + Split-Feature Randomization

2. What is **sample bagging**?
3. What is **split-feature randomization**?
4. How do these techniques affect the bias and variance of an individual tree?
5. How do these techniques affect the bias and variance of an ensemble of trees?

6. For each data point $\mathbf{x}^{(i)}$, define $t^{(-i)}$ to be the set of decision trees that $\mathbf{x}^{(i)}$ was not used to train. Use each tree in $t^{(-i)}$ to make a prediction for $\mathbf{x}^{(i)}$, and use these predictions to make an aggregated prediction $\overline{t^{(-i)}}(\mathbf{x}^{(i)})$ (i.e. for classification take the majority vote). Then, we can define the *out-of-bag* error as follows:

$$E_{OOB} = \frac{1}{N} \sum_{i=1}^N 1 \left(\overline{t^{(-i)}}(\mathbf{x}^{(i)}) \neq y^{(i)} \right)$$

Why can we use E_{OOB} for hyperparameter optimization even though it was calculated using training points we used to learn the decision trees with?

7. **Random Forest Example:** Suppose we train a random forest with two decision trees on the following dataset, using the provided bootstrap samples. Assume that for ties, we predict $Y = 1$.

All	X_0	X_1	X_2	X_3	Y
1	1	0	0	0	1
2	0	0	1	0	1
3	0	0	0	1	1
4	0	0	0	0	0
5	0	1	0	1	1

Sample 1	X_0	X_1	X_2	X_3	Y	Sample 2	X_0	X_1	X_2	X_3	Y
1	1	0	0	0	1	3	0	0	0	1	1
4	0	0	0	0	0	4	0	0	0	0	0
5	0	1	0	1	1	5	0	1	0	1	1

- (a) Suppose we train our first tree on Sample 1 and the split feature randomization chooses $\{X_1, X_2\}$ for the feature candidates at the root. What feature will we split on at the root?
- (b) Suppose we then recurse on the left child (with feature value 0) of the root and split feature randomization chooses $\{X_0, X_2\}$ for the feature indices. What feature will we split on?
- (c) Suppose we train our second tree on Sample 2 and the split feature randomization chooses $\{X_2, X_3\}$ for the feature candidates at the root. What feature will we split on at the root?
- (d) What is the training error of the ensemble?
- (e) What is the out of bag error of the ensemble?

2.2 AdaBoost

2.2.1 AdaBoost Definitions

- T : The number of iterations used to train AdaBoost.
- N : The number of training samples.
- $S = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$: The training samples with binary labels ($y^{(i)} \in \{-1, +1\}$).
- $\omega_t^{(i)}$: The weight assigned to training example i at time t . Note that $\sum_i \omega_t^{(i)} = 1$.
- h_t : The weak learner constructed at time t (a function $X \rightarrow \{-1, +1\}$).
- ϵ_t : The weighted (by ω_t) error of h_t .
- $Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$: The normalization factor for the distribution update at time t .
- $\alpha_t = \frac{1}{2} \ln((1 - \epsilon_t)/\epsilon_t)$: The weight assigned to the learner h_t in the composite hypothesis.
- $H_t(x) = (\sum_{t'=1}^t \alpha_{t'} h_{t'}(x)) / (\sum_{t'=1}^t \alpha_{t'})$: The majority vote of the weak learners, rescaled based on the total weights.
- $g_t(x) = \text{sign}(H_t(x))$: The voting classifier decision function.

2.2.2 AdaBoost Weighting

AdaBoost relies on building an ensemble of weak learners, assigning them weights based on their errors during training.

1. Assume we are in the binary classification setting. What happens to the weight $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$ of classifier h_t if its error $\epsilon_t > 0.5$? Why is this useful?

Note that if we can find weak learners h_t with $\epsilon_t < 0.5$ for all t , training error will decrease exponentially fast in the total number of iterations T .

2. AdaBoost also assigns weights $\omega_t^{(i)}$ for each data point. Explain in broad terms how the weights assigned to examples get updated in each iteration.

