



Lecture 14

Learning Theory

Rishabh Iyer

University of Texas at Dallas

Based on the slides of Vibhav Gogate, David Sontag and Nick Rouzzi

Learning Theory

- So far, we've been focused only on algorithms for finding the best hypothesis in the hypothesis space
 - **Generalization:** How do we know that the learned hypothesis will perform well on the test set? → train & test belongs to same distribution.
 - **Generalization Bounds:** How many samples do we need to make sure that we learn a good hypothesis?

① Given train / test dataset, bound test error
using train error
(# samples)

② Minimum # examples needed for $\leq \epsilon$ gap bet
train - test

Learning Theory

- If the training data is linearly separable, we saw that perceptron/SVMs will always perfectly classify the training data
 - This does not mean that it will perfectly classify the test data
 - Intuitively, if the true distribution of samples is linearly separable, then seeing more data should help us do better

Problem Complexity

- Complexity of a learning problem depends on
 - Size/expressiveness of the hypothesis space (Loosely number of parameters)
 - Accuracy to which a target concept must be approximated
 - Probability with which the learner must produce a successful hypothesis
 - Manner in which training examples are presented, e.g. randomly or by query to an oracle

Goal [Sample Complexity]

① $\text{Err}(\text{test}) \leq \text{Err}(\text{train}) + T(M, C, \delta)$

\uparrow

True-error
with prob atleast
 $1 - \delta$

train data
points

Complexity
of hyp

② $M \leq F(\text{Err}(\text{test}), \text{Err}(\text{train}), \text{comp}, \delta)$ space.

with prob atleast $1 - \delta$.

[Worst Case Bound]

Problem Complexity



- Measures of complexity
 - Sample complexity
 - How much data you need in order to (with high probability) learn a good hypothesis
 - Computational complexity
 - Amount of time and space required to accurately solve (with high probability) the learning problem $\in \text{Poly}(M, d, h)$
 - Higher sample complexity means higher computational complexity

PAC Learning



- Probably approximately correct (PAC)
 - The only reasonable expectation of a learner is that with high probability it learns a close approximation to the **target concept**
 - Specify two small parameters, ϵ and δ , and require that with probability at least $(1 - \delta)$ a system learn a concept with error at most ϵ

Consistent Learners

- Imagine a simple setting
 - The hypothesis space is finite (i.e., $|H| = c$)
 - The true distribution of the data is $p(\vec{x})$, no noisy labels
 - We learned a perfect classifier on the training set, let's call it $h \in H$
 - A learner is said to be **consistent** if it always outputs a perfect classifier (assuming that one exists) $\rightarrow TE = 0$
 - Want to compute the (expected) error of the classifier

We can find a hypothesis with Train error = 0
Can we bound Test error $\leq \epsilon$

Notions of Error

- Training error of $h \in H$
 - The error on the training data
 - Number of samples incorrectly classified divided by the total number of samples
- True error of $h \in H$
 - The error over all possible future random samples
 - Probability, with respect to the data generating distribution, that h misclassifies a random data point

$$(x_1, y_1), \dots, (x_M, y_M)$$

$$\text{Err}_h(\text{Train}) = \frac{\sum_{i=1}^M \text{Error}(x_i, y_i, h)}{M}$$

$$p(h(x) \neq y) = \sum_{n,y} p(n,y) \text{Error}(n, y, h)$$

Learning Theory

- Assume that there exists a hypothesis in H that perfectly classifies all data points and that $|H|$ is finite
$$h \in H : \mathbb{E}_n = 0$$
- The **version space** (set of consistent hypotheses) is said to be ϵ -exhausted if and only if every consistent hypothesis has true error less than ϵ
 - Want enough samples to guarantee that every consistent hypothesis has error at most ϵ

↑
Worst case error of every $\leq \epsilon$
consistent Hypothesis

Learning Theory

- Let $(x^{(1)}, y^{(1)}), \dots, (x^{(M)}, y^{(M)})$ be M labelled data points sampled independently according to p $C_1^h, \dots, C_m^h \in \{0, 1\}$
- Let C_m^h be a random variable that indicates whether or not the m^{th} data point is correctly classified
- The probability that h misclassifies the m^{th} data point is

$$p(C_m^h = 0) = \sum_{(x,y)} p(x, y) \underbrace{\mathbb{1}_{h(x) \neq y}}_{=} = (\epsilon_h)$$

Learning Theory

J Pave, $n, y \sim P$



- Let $(x^{(1)}, y^{(1)}), \dots, (x^{(M)}, y^{(M)})$ be M labelled data points sampled independently according to p
- Let C_m^h be a random variable that indicates whether or not the m^{th} data point is correctly classified
- The probability that h misclassifies the m^{th} data point is

$$p(C_m^h = 0) = \sum_{(x,y) \sim P} p(x, y) \mathbf{1}_{h(x) \neq y} = \epsilon_h$$

Sample "test points" $\downarrow N$

Compute test error ϵ_n

These points
 $\approx p(C_m^h = 0)$

Probability that a randomly sampled pair (x,y) is incorrectly classified by h

$$E[\mathbf{1}_{h(x) \neq y}]$$

Learning Theory

- Let $(x^{(1)}, y^{(1)}), \dots, (x^{(M)}, y^{(M)})$ be M labelled data points sampled independently according to p
- Let C_m^h be a random variable that indicates whether or not the m^{th} data point is correctly classified
- The probability that h misclassifies the m^{th} data point is

$$p(C_m^h = 0) = \sum_{(x,y)} p(x,y) \mathbf{1}_{h(x) \neq y} = \epsilon_h$$

This is the true error of hypothesis h

Learning Theory

- Probability that all data points classified correctly?
- Probability that a hypothesis $h \in H$ whose true error is at least ϵ correctly classifies the m data points is then

Learning Theory

- Probability that all data points classified correctly?

$$p(C_1^h = 1, \dots, C_M^h = 1) = \underbrace{\prod_{m=1}^M p(C_m^h = 1)}_{h \text{ is consistent}} = (1 - \epsilon_h)^M$$

- Probability that a hypothesis $h \in H$ whose true error is at least ϵ correctly classifies the m data points is then

Learning Theory

- Probability that all data points classified correctly?

$$p(C_1^h = 1, \dots, C_M^h = 1) = \prod_{m=1}^M p(C_m^h = 1) = (1 - \epsilon_h)^M$$

- Probability that a hypothesis $h \in H$ whose true error is at least ϵ correctly classifies the m data points is then $p(C_i^h = 0) \geq \underline{\epsilon}$

$$p(C_1^h = 1, \dots, C_M^h = 1) \leq \underbrace{(1 - \epsilon)^M}_{\leq e^{-\epsilon M}}$$

for $\epsilon \leq 1$

$$1 - \epsilon \leq e^{-\epsilon}$$

$$e^n \approx 1 + n$$

The Union Bound

- Let $H_{BAD} \subseteq H$ be the set of all hypotheses that have true error at least ϵ
 $\forall h \in H_{BAD}, \quad \epsilon_h \geq \epsilon$
- From before for each $h \in H_{BAD}$,
 $p(h \text{ correctly classifies all } M \text{ data points}) \leq e^{-\epsilon M}$
- So, the probability that $\underline{\text{some}} \ h \in H_{BAD}$ correctly classifies all of the data points is

$$\begin{aligned}
 p\left(\bigvee_{h \in H_{BAD}} (C_1^h = 1, \dots, C_M^h = 1)\right) &\leq \sum_{h \in H_{BAD}} p(C_1^h = 1, \dots, C_M^h = 1) \\
 &\leq |H_{BAD}| e^{-\epsilon M} \\
 &\leq |H| e^{-\epsilon M}
 \end{aligned}$$

$\underbrace{\quad}_{Eh.}$

Union Bound.

Set of Events $h \in H$.

$$P\left(\bigvee_{h \in H} E_h\right) \leq \sum_{h \in H} P(E_h)$$

- What we just proved:
 - **Theorem:** For a finite hypothesis space, \underline{H} , with M i.i.d. samples, and $0 < \epsilon < 1$, the probability that the version space is not ϵ -exhausted is at most $|H|e^{-\epsilon M}$
 - We can turn this into a **sample complexity bound**

- What we just proved:

- **Theorem:** For a finite hypothesis space, \underline{H} , with M i.i.d. samples, and $0 < \epsilon < 1$, the probability that **there exists a hypothesis in H that is consistent with the data but has true error larger than ϵ** is at most $|H|e^{-\epsilon M} \leq \delta$

$$|H|e^{-\epsilon M} \leq \delta$$

↑ bounded ↑ Trainerror(h) = 0
 ↓ ↓
 Finite. M training points

- We can turn this into a **sample complexity bound**

$H_{\text{BAD}} \rightarrow \text{BAD}$ from generalization perspective.

$$l_h \geq \epsilon$$

$$|H| e^{-\varepsilon M} \leq \delta$$

$$\log |H| - \varepsilon M \leq \log \delta$$

$$\log |H| - \log \delta \leq \varepsilon M$$

$$\Rightarrow M \geq \frac{1}{2} [\log |H| - \log \delta]$$

$$\geq \frac{1}{2} \left[\log |H| + \log \frac{1}{\delta} \right]$$

Sample Complexity (Bound on # train examples)



- Let δ be an upper bound on the desired probability of not ϵ -exhausting the sample space
 - That is, the probability that the version space is not ϵ -exhausted is at most $|H|e^{-\epsilon M} \leq \delta$
- Solving for M yields

$$\underline{M} \geq -\frac{1}{\epsilon} \ln \frac{\delta}{|H|}$$

$$= \left(\ln |H| + \ln \frac{1}{\delta} \right) / \epsilon$$

$$\begin{aligned} \epsilon &= 10^3 \\ \delta &= 10^{-6} \\ |H| &= 10^{10} \end{aligned}$$

$$M \geq (10 + 6) \times 10^3$$

$1 - \delta = \text{Prob that}$
 $\text{every } h \in H$
 $\text{that is consistent}$
 has true error
 $\leq \epsilon$

Sample Complexity

- Let δ be an upper bound on the desired probability of not ϵ -exhausting the sample space
 - That is, the probability that the version space is not ϵ -exhausted is at most $\underbrace{|H|}_{} e^{-\epsilon M} \leq \delta$
- Solving for M yields

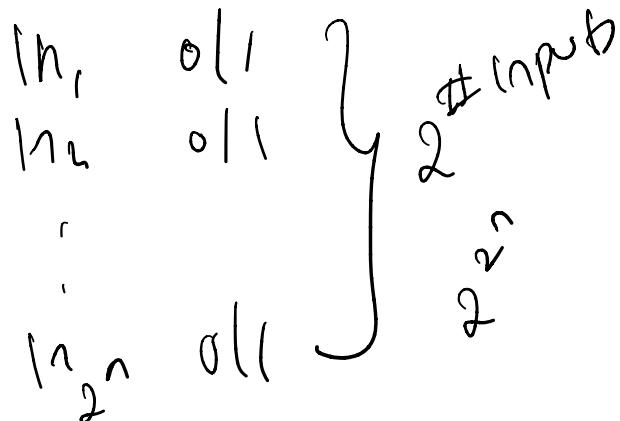
$$\begin{aligned} M &\geq -\frac{1}{\epsilon} \ln \frac{\delta}{|H|} \\ &= \left(\ln |H| + \ln \frac{1}{\delta} \right) / \epsilon \end{aligned}$$

Typical in learning
theory

This is sufficient,
but not necessary
(union bound is
quite loose)

Decision Trees

- Suppose that we want to learn an arbitrary Boolean function given n Boolean features
- Hypothesis space consists of all decision trees
 - Size of this space = ?
 - How many samples are sufficient?



$$\begin{aligned}
 & n_1 \ n_2 \dots \ n_n \quad y \\
 & 0/1 \quad 0/1 \quad \dots \quad 0/1 \quad 0/1 \\
 & \text{\# possible input values} \\
 & n_1 \cdot n_2 \cdots n_n = 2^n \\
 & \text{\# Hypothesis functions} = 2^{2^n}
 \end{aligned}$$

Decision Trees

- Suppose that we want to learn an arbitrary Boolean function given n Boolean features
- Hypothesis space consists of all decision trees
 - Size of this space = $\frac{2^{2^n}}{\gamma}$ = number of Boolean functions on n inputs
Doubly exponential
- How many samples are sufficient?

$$M \geq \left(\ln(2^{2^n}) + \ln \frac{1}{\delta} \right) / \epsilon \quad \leftarrow$$

$$M \geq \left(2^n + \ln \frac{1}{\delta} \right) / \epsilon$$

Generalizations

train-error = 0

- How do we handle situations with no perfect classifier?
 - Pick the hypothesis with the lowest error on the training set
- What do we do if the hypothesis space isn't finite?
 - Infinite sample complexity? $|H| = \infty$
 - Coming soon...

- ① DT on red dots
- ② Perceptron / SVM / LR / Lin Reg

Chernoff Bounds

- Chernoff bound: Suppose Y_1, \dots, Y_M are i.i.d. random variables taking values in $\{0, 1\}$ such that $E_p[Y_i] = y$. For $\epsilon > 0$,

$$p \left(\left| y - \frac{1}{M} \sum_{m=1}^M Y_m \right| \geq \epsilon \right) \leq 2e^{-2M\epsilon^2}$$

Chernoff Bounds

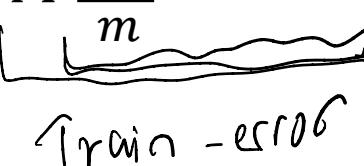
- Chernoff bound: Suppose Y_1, \dots, Y_M are i.i.d. random variables taking values in $\{0, 1\}$ such that $E_p[Y_i] = y$. For $\epsilon > 0$,

$$p\left(\left|y - \frac{1}{M} \sum_m Y_m\right| \geq \epsilon\right) \leq 2e^{-2M\epsilon^2}$$

$Y_i = 1 - C_i^h$
 ↑
 error of item i

- Applying this to $1 - C_1^h, \dots, 1 - C_M^h$ gives

$$p\left(\left|\epsilon_h - \frac{1}{M} \sum_m (1 - C_m^h)\right| \geq \epsilon\right) \leq 2e^{-2M\epsilon^2}$$



Chernoff Bounds

- Chernoff bound: Suppose Y_1, \dots, Y_M are i.i.d. random variables taking values in $\{0, 1\}$ such that $E_p[Y_i] = y$. For $\epsilon > 0$,

$$p\left(\left|y - \frac{1}{M} \sum_m Y_m\right| \geq \epsilon\right) \leq 2e^{-2M\epsilon^2}$$

- Applying this to $1 - C_1^h, \dots, 1 - C_M^h$ gives

$$p\left(\epsilon_h - \frac{1}{M} \sum_m (1 - C_m^h) \geq \epsilon\right) \leq e^{-2M\epsilon^2}$$

↙ ↘
Hoeffding Bound

This is the training error

Test error

↑
train - test

PAC Bounds [Sample Complexity]



- **Theorem:** For a finite hypothesis space H finite, M i.i.d. samples, and $0 < \epsilon < 1$, the probability that true error of any of the best classifiers (i.e., lowest training error) is larger than its training error plus ϵ is at most $|H|e^{-2M\epsilon^2}$
- Sample complexity (for desired $\delta \geq |H|e^{-2M\epsilon^2}$)

$$M \geq \left(\ln|H| + \ln \frac{1}{\delta} \right) / 2\epsilon^2$$

$\vdash \delta \Rightarrow \text{Prob that } \text{test } E_n - \text{train } E_n \leq \epsilon$

$$\epsilon \geq \sqrt{\frac{1}{M} (\log|H| + \log \frac{1}{\delta})}$$

PAC Bounds [Generalization Bound]



- If we require that the previous error is bounded above by δ , then with probability $(1 - \delta)$, for all $h \in H$

$$\text{true-} \epsilon_h \leq \underbrace{\epsilon_h^{\text{train}}}_{\text{"bias"}} + \sqrt{\underbrace{\frac{1}{2M} \left(\ln |H| + \ln \frac{1}{\delta} \right)}_{\text{"variance"} \text{ [Complexity of } H]}}$$

“variance” [Complexity of H]

- For small $|H|$
 - High bias (may not be enough hypotheses to choose from)
 - Low variance

$$\epsilon_h = \epsilon_h^{\text{train}} + \underline{\epsilon}$$

Under fitting

PAC Bounds

- If we require that the previous error is bounded above by δ , then with probability $(1 - \delta)$, for all $h \in H$

$$\epsilon_h \leq \epsilon_h^{train} + \sqrt{\frac{1}{2M} \left(\ln |H| + \ln \frac{1}{\delta} \right)}$$

A diagram illustrating the decomposition of the PAC bound. The term ϵ_h^{train} is bracketed on the left and labeled "bias". The term $\sqrt{\frac{1}{2M} \left(\ln |H| + \ln \frac{1}{\delta} \right)}$ is bracketed on the right and labeled "variance".

- For large $|H|$ [More Complex Hypotheses space]
 - Low bias (lots of good hypotheses)
 - High variance

↑ overfitting

PAC Bounds



- If we require that the previous error is bounded above by δ , then with probability $(1 - \delta)$, for all $h \in H$

$$\epsilon_h \leq \epsilon_h^{train} + \sqrt{\frac{1}{2M} \left(\ln |H| + \ln \frac{1}{\delta} \right)}$$

“bias”

“variance”

- For large $|H|$
 - Low bias (lots of good hypotheses)
 - High variance

$M = \text{small}$
↓
 $\text{train error} = \text{low}$
 $\text{Var} = \text{high}$

$M = \text{large}$
↓
 $\text{train error} = \text{high}$
 $\text{Var} = \text{low}$