Learning Theory

Aarti Singh

Machine Learning 10-315 Nov 22, 2021

Slides courtesy: Carlos Guestrin





Learning Theory

- We have explored many ways of learning from data
- But...
 - Can we certify how good is our classifier, really?
 - How much data do I need to make it "good enough"?

PAC Learnability

Probably Approximately Correct

- True function space, F fer
- Model space, H LeH

F is PAC Learnable by a learner using H if

there exists a learning algorithm s.t. for all functions in

 \digamma \not H, for all distributions over inputs, for all $0 < \varepsilon$, $\delta < 1$,

with probability > $1-\delta$, the algorithm outputs a model

 $h \in H \text{ s.t. error}_{true}(h) \leq \varepsilon$

in time and samples that are polynomial in $1/\epsilon$, $1/\delta$ and n.

A simple setting

- Classification
 - m i.i.d. data points
 - Finite number of possible classifiers in model class (e.g., dec. trees of depth d)
- Lets consider that a learner finds a classifier h that gets zero error in training
 - $-\operatorname{error}_{\operatorname{train}}(h) = 0$

- What is the probability that h has more than ε true (= test) error?
 - $error_{true}(h) ≥ ε$

3 < (r + (x) 1) > E

How likely is a bad classifier to get m data points right?

- P(h(x)+Y) ≥ ε• Consider a bad classifier h i.e. error_{true}(h) ≥ ε
- Probability that h gets one data point right

Probability that h gets m data points right

How likely is a learner to pick a bad classifier?

Usually there are many (say k) bad classifiers in model class

$$h_1, h_2, ..., h_k$$
 s.t. $error_{true}(h_i) \ge \varepsilon$ $i = 1, ..., k$

 Probability that learner picks a bad classifier = Probability that some bad classifier gets 0 training error

```
Prob(h_1 gets 0 training error OR h_2 gets 0 training error OR ... OR h_k gets 0 training error)
```

≤ Prob(
$$h_1$$
 gets 0 training error) +

Prob(h_2 gets 0 training error) + ... +

Prob(h_k gets 0 training error) $\angle k(1 - \epsilon)^m$

P(AUR)
P(A)+P(R)

Union bound Loose but works

How likely is a learner to pick a bad classifier?

Usually there are many many (say k) bad classifiers in the class

$$h_1, h_2, ..., h_k$$

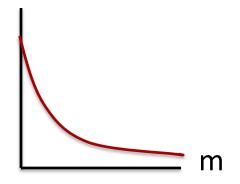
s.t.
$$error_{true}(h_i) \ge \varepsilon$$
 $i = 1, ..., k$

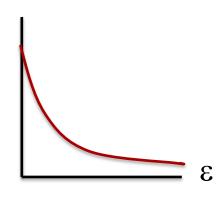
Probability that learner picks a bad classifier

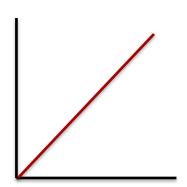
1-850-8

$$\leq k (1-\epsilon)^m \leq |H| (1-\epsilon)^m \leq |H| e^{-\epsilon m}$$

→ Size of model class







PAC (Probably Approximately Correct) bound

• Theorem [Haussler'88]: Model class H finite, dataset D with m i.i.d. samples, $0 < \varepsilon < 1$: for any learned classifier h that gets 0 training error:

$$P(\text{error}_{true}(h) \ge \epsilon) \le |H|e^{-m\epsilon} \le \delta$$

• Equivalently, with probability $\ \geq 1-\delta$

Important: PAC bound holds for all h with 0 training error, but doesn't guarantee that algorithm finds best h!!!

Using a PAC bound

$$|H|e^{-m\epsilon} \le \delta$$

• Given ε and δ , yields sample complexity

#training data,
$$m \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{\epsilon}$$

• Given m and δ , yields error bound

error,
$$\epsilon \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$$

Poll

Ghartion, (N/=0

Assume m is the minimum number of training examples sufficient to guarantee that with probability $1-\delta$ a consistent learner using model class H will output a classifier with true error at worst ϵ .

Then a second learner that uses model space H' will require 2m training examples (to make the same guarantee) if |H'| = 2|H|.

A. True

B. False

If we double the number of training examples to 2m, the error bound ϵ will be halved.

C. True

D. False

Limitations of Haussler's bound

Only consider classifiers with 0 training error

h such that zero error in training, $error_{train}(h) = 0$

Dependence on size of model class |H|

$$m \ge \frac{\ln|H| + \ln\frac{1}{\delta}}{\epsilon}$$

what if |H| too big or H is continuous (e.g. linear classifiers)?

PAC bounds for finite model classes

H - Finite model class

e.g. decision trees of depth k histogram classifiers with binwidth h

With probability $\geq 1-\delta$,

1) For all
$$h \in H$$
 s.t. $error_{train}(h) = 0$,

error_{true}(h)
$$\leq \varepsilon = \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$$

Haussler's bound

What if our classifier does not have zero error on the training data?

- A learner with zero training errors may make mistakes in test set
- What about a learner with error_{train}(h) ≠ 0 in training set?
- The error of a classifier is like estimating the parameter of a coin!

$$error_{true}(h) := P(h(X) \neq Y) \equiv P(H=1) =: \theta$$
 $error_{train}(h) := \frac{1}{m} \sum_{i} \mathbf{1}_{h(X_i) \neq Y_i} \equiv \frac{1}{m} \sum_{i} Z_i =: \widehat{\theta}$

Hoeffding's bound for a single classifier

• Consider m i.i.d. flips $x_1,...,x_m$, where $x_i \in \{0,1\}$ of a coin with parameter θ . For $0 < \epsilon < 1$:

$$P\left(\left|\theta - \frac{1}{m}\sum_{i}x_{i}\right| \ge \epsilon\right) \le 2e^{-2m\epsilon^{2}}$$

• Central limit theorem:
$$Z_{(1,-2m)} \stackrel{\text{iid}}{=} E[2i] = \mu \text{ vol}(z_i) = \sigma^2$$

$$\lim_{m \to \infty} (\frac{1}{m} \stackrel{\text{iid}}{=} z_i - \mu) \longrightarrow \mathcal{N}(0, \sigma^2) = \lim_{m \to \infty} \frac{1}{2} z_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

Rev. ii $\mu = 0$ $\sigma^2 = 0(1-0) \leq \frac{1}{4}$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

$$\lim_{n \to \infty} \frac{1}{m} \stackrel{\text{iid}}{=} x_i \longrightarrow \mathcal{N}(\mu, \sigma_m)$$

Hoeffding's bound for a single classifier

• Consider m i.i.d. flips $x_1,...,x_m$, where $x_i \in \{0,1\}$ of a coin with parameter θ . For $0 < \epsilon < 1$:

$$P\left(\left|\theta - \frac{1}{m}\sum_{i}x_{i}\right| \geq \epsilon\right) \leq 2e^{-2m\epsilon^{2}}$$

$$\theta = P(h(x) \neq Y) = e^{-2m\epsilon^{2}}$$

For a single classifier h

$$P\left(|\operatorname{error}_{true}(h) - \operatorname{error}_{train}(h)| \ge \epsilon\right) \le 2e^{-2m\epsilon^2}$$

Hoeffding's bound for |H| classifiers

• For each classifier h_i:

$$P\left(\left|\operatorname{error}_{true}(h_i) - \operatorname{error}_{train}(h_i)\right| \ge \epsilon\right) \le 2e^{-2m\epsilon^2}$$

- What if we are comparing |H| classifiers?
 Union bound
- *Theorem*: Model class H finite, dataset D with m i.i.d. samples, $0 < \varepsilon < 1$: for any learned classifier $h \in H$:

$$P\left(\operatorname{error}_{true}(h) - \operatorname{error}_{train}(h)| \ge \epsilon\right) \le 2|H|e^{-2m\epsilon^2} \le \delta$$

Important: PAC bound holds for all h, but doesn't guarantee that $_{16}$ algorithm finds best h!!!

Summary of PAC bounds for finite model classes

With probability $\geq 1-\delta$,

1) For all $h \in H$ s.t. $error_{train}(h) = 0$,

error_{true}(h)
$$\leq \varepsilon = \frac{\ln|H| + \ln\frac{1}{\delta}}{m}$$

Haussler's bound

2) For all $h \in H$ $|error_{true}(h) - error_{train}(h)| \le \varepsilon = \sqrt{\frac{\ln|H| + \ln\frac{2}{\delta}}{2m}}$

Hoeffding's bound

PAC bound and Bias-Variance tradeoff

$$P\left(|\operatorname{error}_{true}(h) - \operatorname{error}_{train}(h)| \ge \epsilon\right) \le 2|H|e^{-2m\epsilon^2} \le \delta$$

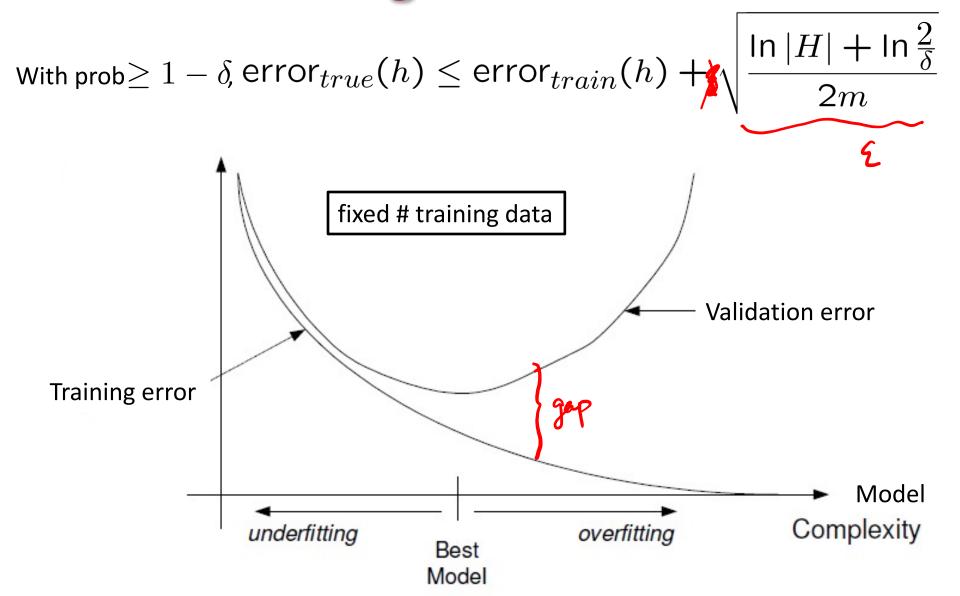
• Equivalently, with probability $\geq 1-\delta$

$$\operatorname{error}_{true}(h) \leq \operatorname{error}_{train}(h) + \sqrt{\frac{\ln|H| + \ln\frac{2}{\delta}}{2m}}$$
 ed m

Fixed m

Model class	↓	↓
complex	small	large
simple	large	small

Training vs. Test Error



What about the size of the model class?

$$2|H|e^{-2m\epsilon^2} \le \delta$$

Sample complexity

$$m \ge \frac{1}{2\epsilon^2} \left(\ln|H| + \ln\frac{2}{\delta} \right)$$

How large is the model class?

Number of decision trees of depth k

Recursive solution:
Given *n* binary attributes

$$m \ge \frac{1}{2\epsilon^2} \left(\ln|H| + \ln\frac{2}{\delta} \right)$$

 H_k = Number of **binary** decision trees of depth k

$$H_0^- = 2$$

 $H_k = (\text{\#choices of root attribute})$

- *(# possible left subtrees)
- *(# possible right subtrees) = $n * H_{k-1} * H_{k-1}$

Write
$$L_k = log_2 H_k$$

$$L_0 = 1 = \log_2 2$$

$$L_k = \log_2 n + 2L_{k-1} = \log_2 n + 2(\log_2 n + 2L_{k-2})$$

$$= \log_2 n + 2\log_2 n + 2^2\log_2 n + ... + 2^{k-1}(\log_2 n + 2L_0)$$

So
$$L_k = (2^k-1)(1+\log_2 n) + 1$$

PAC bound for decision trees of depth k

$$m \ge \frac{\ln 2}{2\epsilon^2} \left((2^k - 1)(1 + \log_2 n) + 1 + \log_2 \frac{2}{\delta} \right)$$

- Bad!!!
 - Number of points is exponential in depth k!

But, for m data points, decision tree can't get too big...

Number of leaves never more than number data points

Number of decision trees with k leaves

$$\boxed{m} \ge \frac{1}{2\epsilon^2} \left(\ln|H| + \ln\frac{2}{\delta} \right)$$

 H_k = Number of binary decision trees with k leaves

$$H_1 = 2$$

 H_k = (#choices of root attribute) *

[(# left subtrees wth 1 leaf)*(# right subtrees wth k-1 leaves)

- + (# left subtrees wth 2 leaves)*(# right subtrees wth k-2 leaves)
- + ...
- + (# left subtrees wth k-1 leaves)*(# right subtrees wth 1 leaf)]

$$H_k = n \sum_{i=1}^{k-1} H_i H_{k-i} = n^{k-1} C_{k-1}$$
 (C_{k-1}: Catalan Number)

Loose bound (using Sterling's approximation):

$$H_k \leq n^{k-1} 2^{2k-1} \qquad |\mathsf{n}| \mathsf{Hu}| \leq (\mathsf{k}-\mathsf{i}) \, \mathsf{ln} \, \mathsf{n}$$

Number of decision trees

With k leaves

$$m \ge \frac{1}{2\epsilon^2} \left(\ln|H| + \ln\frac{2}{\delta} \right)$$

$$\log_2 H_k \le (k-1)\log_2 n + 2k - 1$$
 linear in k number of points m is linear in #leaves

With depth k

$$log_2 H_k = (2^k-1)(1+log_2 n) +1$$
 exponential in k number of points m is exponential in depth

What did we learn from decision trees?

Moral of the story:

Complexity of learning not measured in terms of size of model space, but in maximum <u>number of points</u> that allows consistent classification

Rademacher Complexity

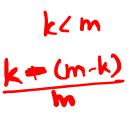
 Instead of all possible labelings, measure complexity by how accurately a model space can match a random labeling of the data.

For each data point i, draw random label

$$\sigma_i$$
 s.t. $P(\sigma_i = +1) = \frac{1}{2} = P(\sigma_i = -1)$

Then empirical Rademacher complexity of H is

$$\widehat{R}_{m}(H) = \mathbb{E}_{\sigma} \left[\sup_{h \in H} \left(\frac{1}{m} \sum_{i=1}^{m} \sigma_{i} h(X_{i}) \right) \right]$$



Finite model class

 Rademacher complexity can be upper bounded in terms of model class size |H|:

$$\widehat{R}_m(H) \le \sqrt{\frac{2\ln|H|}{m}}$$

Often Rademacher bounds are significantly better