

Bias and VC dimension

Machine Learning Theory (MLT)

Edinburgh

Rik Sarkar

Today's topics

- Uniform convergence
- No free lunch theorem
- Error and bias-complexity tradeoff
- VC dimension and fundamental theorem of statistical learning

Uniform convergence

- S is ϵ –representative w.r.t $(\mathcal{Z}, \mathcal{H}, \mathcal{D})$ if:
 - $\forall h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$
- S gives a good estimate of the true loss for each h
- \mathcal{H} has uniform convergence if there is $m_{\mathcal{H}}^{UC}: (0,1)^2 \rightarrow \mathbb{N}$
 - Such that a random sample $S \sim \mathcal{D}^m$ of size $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$
 - Is ϵ –representative with probability at least $1 - \delta$
- A hypothesis class is said to have uniform convergence
 - if for any ϵ, δ :
 - Any random sample S of large enough sample size m (a function of ϵ, δ)
 - Is ϵ – representative with probability of at least $1 - \delta$
 - Will give a good estimate of

Corollary

- If \mathcal{H} has uniform convergence with $m_{\mathcal{H}}^{UC}$,
 - Then \mathcal{H} is PAC learnable with $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)$
- So, every finite \mathcal{H} is learnable
 - In the sense that based on S
 - ERM algorithm finds an $h \in \mathcal{H}$
 - Whose true loss is close to the best general $h^* \in \mathcal{H}$

- Theorem:
- Every finite \mathcal{H} has uniform convergence
 - i.e. Given a random S , $\mathbb{P}[\exists h \in \mathcal{H}: |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon] \leq \delta$
- (And therefore every finite \mathcal{H} is agnostic PAC-learnable)
- To prove this, we need the Chernoff-hoeffding bound

Chernoff-Hoeffding bound

- Very important result in theoretical CS and ML
- Suppose θ_i are random variables with average $\frac{1}{m} \sum_{i=1}^m \theta_i$
- Suppose μ is the expected value
- Law of large numbers: with increasing m , $\frac{1}{m} \sum_{i=1}^m \theta_i$ approaches μ
 - I.e., $\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right|$ becomes smaller
- But how fast? What m do we need to get ϵ -close to μ ?
- Chernoff-Hoeffding bound:
 - $\mathbb{P} \left[\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon \right] \leq 2e^{-2m\epsilon^2}$

- Proof that $\mathbb{P}[\exists h \in \mathcal{H}: |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon] \leq \delta$
- Take any $h \in \mathcal{H}$
- Now take a random sample S
- Let us write $\mu = \mathbb{E}[L_S(h)] = L_{\mathcal{D}}(h)$
 - I.e. note that the expected value of empirical loss is the true loss
- For every $z_i \in S$, we write its loss on h as θ_i . I.e. $\theta_i = \ell(h, z_i)$
- Then the empirical loss is $L_S(h) = \frac{1}{m} \sum_{i=1}^m \theta_i$
- So, what is the probability that $\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon$?

- What is the probability that $\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon$?
- Using Chernoff bound, probability that any one h has large error is:
 - $\mathbb{P} \left[\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon \right] \leq 2e^{-2m\epsilon^2}$
- Summing over all $h \in \mathcal{H}$, probability that one or more has large error is:

- What is the probability that $\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon$?
- Using Chernoff bound, probability that any one h has large error is:
 - $\mathbb{P} \left[\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon \right] \leq 2e^{-2m\epsilon^2}$
- Summing over all $h \in \mathcal{H}$, probability that one or more has large error is:
 - $\leq 2|\mathcal{H}|e^{-2m\epsilon^2}$ (by union bound)
- Substitute $m \geq \frac{1}{2\epsilon^2} \log \left(\frac{2|\mathcal{H}|}{\delta} \right)$ to get a probability bound δ

- So, we have proved finite classes are all PAC learnable
- How useful is it?

- So, we have proved finite classes are all PAC learnable
- How useful is it?

- \mathcal{H} is described by some parameters
 - E.g. coefficients of polynomials, weights on edges...
- We can always discretize by taking closely spaced discrete points
- In fact in a computer, numbers are always discretized
 - And thus practical hypothesis classes are technically finite
- Though searching the entire class may not be practical...

Do we need to think about hypothesis classes?

- In class, we discussed why we cannot have “all possible” hypotheses as our class.
 - Challenges in computation
 - Challenges in enumerating
- We now see a theorem that says there is no universal PAC learner
- For any learner A , there is an input on which it fails to find a good model

No Free Lunch theorem (thm 5.1 in book)

- For 0-1 loss. Assume $m < \frac{|\mathcal{X}|}{2}$
- There is \mathcal{D}, f such that
 - $L_{\mathcal{D}}(f) = 0$ (i.e. f is a perfect classifier that exists)
 - And we would like to find something close to f
 - For a random S
 - With probability $\geq \frac{1}{7}$, the true loss $L_{\mathcal{D}}(A(S)) \geq \frac{1}{8}$
- This violates the ϵ, δ guarantees of PAC
- We skip the proof

Conclusion: Prior knowledge is necessary

- \mathcal{H} represents what we know or can guess about the problem
- A restricted class is suitable when we have a pretty good idea, or prior knowledge
- A larger class is suitable when we have less knowledge about the problem
- Taking the set of all possible hypotheses or functions will imply no knowledge of domain
- Corollary (5.2 in book): If \mathcal{X} is infinite, and \mathcal{H} is set of all possible functions from \mathcal{X} to $\{0,1\}$, then \mathcal{H} is not PAC learnable

Impact of Prior knowledge

- So, we have to choose a fixed \mathcal{H}
 - We make some assumptions about the type of solution that can work
 - This introduces a *Bias* we are looking for certain types of solutions instead of all possible solutions

Error decomposition

- True loss: $L_{\mathcal{D}}(h_S) = \epsilon_{app} + \epsilon_{est}$
- Approximation error $\epsilon_{app} = L_{\mathcal{D}}(h^*)$
 - Min true error in the hypothesis class
 - Limitation of the choice of hypothesis class
- Estimation error $\epsilon_{est} = L_{\mathcal{D}}(h_S) - L_{\mathcal{D}}(h^*)$:
 - Difference between approximation error and true error
 - Error due to sampling and overfitting choosing suboptimal h_S

Bias complexity tradeoff

- Rich/complex \mathcal{H}
 - Small approximation error
 - Large estimation error (due to overfitting)
 - Needs more data
- Small, restricted \mathcal{H}
 - Small estimation error
 - Needs less data
 - Large approximation error
- A lot of ML is about designing good \mathcal{H} , balancing errors
 - Often by making use of our knowledge of the domain

Infinite hypothesis classes

- We have proved that finite \mathcal{H} are PAC learnable
- What about infinite \mathcal{H} ?

- We already proved PAC guarantee for an infinite \mathcal{H}
 - The threshold classifier for ripe papayas: real numbers $[t, 1]$
- Why were we able to prove the guarantees?

Observation

- Sample complexity $m \geq \frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\epsilon^2}$ can be written as
 - $m \geq \frac{\log(2|\mathcal{H}|) + \log\left(\frac{1}{\delta}\right)}{2\epsilon^2}$
- It has two components
 - Complexity of hypothesis class
 - Confidence probability

- Dimensionality
 - A measure of complexity of \mathcal{H}
 - Allows us to get efficient results in the 1-D case
- Another measure of complexity was $|\mathcal{H}|$ -- for finite classes
- VC-dimension: A complexity measure for infinite classes
 - Ability of \mathcal{H} to split different arrangements of points into different subsets

Shattering

- Take a point set $C \subset \mathcal{X}$
- C is shattered by \mathcal{H} if
- Any classification of points in C can be achieved by \mathcal{H}
- That is, for each possible 0-1 labelling of points in C
 - There is an $h \in \mathcal{H}$ that selects all of the ones and none of the zeros

VC dimension

- VC dimension of \mathcal{H} is
- The size of the largest set $C \subset \mathcal{X}$ that can be shattered by \mathcal{H}
- For VC dim to be d , we have to show:
 - There is one set of size d that is shattered by \mathcal{H}
 - No set of size $d + 1$ is shattered by \mathcal{H}

VC dim examples

- Threshold functions: Dim 1
- Intervals : Dim 2
- Axis aligned rectangles: Dim 4

Finite classes

- On C there are $2^{|C|}$ possibly binary classifications
- Thus, C cannot be shattered if $|\mathcal{H}| < 2^{|C|}$
- Therefore: $VCdim(\mathcal{H}) \leq \log_2 |\mathcal{H}|$

Number of parameters

- Number of parameters of \mathcal{H} is a good measure of complexity
- Often equals VCdim
 - But not always

Fundamental theorem of statistical learning

THEOREM 6.7 (The Fundamental Theorem of Statistical Learning) *Let \mathcal{H} be a hypothesis class of functions from a domain \mathcal{X} to $\{0, 1\}$ and let the loss function be the 0 – 1 loss. Then, the following are equivalent:*

- 1. \mathcal{H} has the uniform convergence property.*
- 2. Any ERM rule is a successful agnostic PAC learner for \mathcal{H} .*
- 3. \mathcal{H} is agnostic PAC learnable.*
- 4. \mathcal{H} is PAC learnable.*
- 5. Any ERM rule is a successful PAC learner for \mathcal{H} .*
- 6. \mathcal{H} has a finite VC-dimension.*

In more detail

1. \mathcal{H} has the uniform convergence property with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}^{uc}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

2. \mathcal{H} is agnostic PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

3. \mathcal{H} is PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$$

- Proof: Omitted
- Other types of loss functions:
 - Similar properties hold. But do search for exact results before use.

Structural Risk Minimization

- We have till now treated \mathcal{H} as binary choice: hypothesis in the class are all equally valid, while outside the class are disallowed
- What if all $h \in \mathcal{H}$ are not equally desirable?
- Structural Risk Minimization
 - Assign different preferences to different hypothesis
- Examples:
 - Assign a weight $w(h)$ to each hypothesis. Higher weight reflects higher preference
 - Divide hypothesis class into subclasses, assign a weight to each class

Some example weighting functions

- Polynomial degree
 - We would usually prefer lower degree polynomials
- Polynomial coefficients
 - We prefer smaller coefficients as that describes gentler/smoother functions
- Minimum description length
 - How many bits or characters does it take to represent the class?
 - Shorter description length (simpler class) is more desirable

Regularization

- A common way to achieve simpler models
- Include a “penalty” function to the loss
- E.g. sum of coefficients, or description length...

- Minimizing the loss now includes minimizing the actual loss and minimizing the penalty

- More on regularization later

Discussion: PAC, VCDim and ML in practice

- The issue: PAC and VC analysis does not work too well on Deep learning
- VC dim of neural networks are hard to compute. A simple bound is $VCDim = O(|E|)$
 - Not useful
 - Varies somewhat by activation function etc.
- However, we are studying this because:
 - Knowing the background helps us in building more comprehensive theories of future ML
 - The set up with formal definition of ML is still valid
 - Rigorous introduction to important concepts like “Probably approximately correct”, hypothesis classes, their complexity, impossibility results etc.

Recap till now

- Hypothesis classes
- Empirical and true loss. Empirical loss minimization
- Sample complexity
- PAC learnability (realizable, finite)
- Agnostic PAC learnability (finite)
- Bias- complexity tradeoff
- VC dim
- PAC learning infinite classes: Fundamental theorem of statistical learning
- Next: Algorithms, convex learning, stochastic gradient descent, neural networks