

Learning from Samples

Wednesday, March 22, 2017 8:50 AM

▷ Last time we saw:

- Vickrey with a different reserve for each player is a 2-approx. to Myerson in non i.i.d. regular environments.

- Vickrey with a single reserve $r \sim F$ is a $(1 - \frac{1}{n+1})$ -approx for i.i.d regular where each $v_i \sim F$.

▷ So if we are given a single sample r from the value distribution then we can construct a mechanism M_r s.t.

$$E_{r, \vec{v} \sim F^n} [\text{Rev}_{M_r}(\vec{v})] \geq \frac{1}{2} \underbrace{\sup_M E_{\vec{v} \sim F^n} [\text{Rev}_M(\vec{v})]}_{\text{OPT}(F)}$$
$$\text{OPT}(F) = E_{\vec{v} \sim F^n} [\text{Rev}_{\text{Myerson}}(\vec{v})]$$

▷ What if we had more samples?

$$S = \{ r_1, \dots, r_m \} \sim F^m$$

Can we give a mechanism M_S

such that:

$$E_{S, \vec{v}} [\text{Rev}_{M_S}(\vec{v})] \geq \text{OPT}(D) - \epsilon(m)$$

with $\epsilon(m) \rightarrow 0$?

▷ What about w.h.p. i.e. w.p. $1 - \delta$

over S :

$$E_{\vec{v}} [Rev_{M_S}(\vec{v})] \geq OPT(D) - \epsilon(m, \delta)$$

with $\epsilon(m, \delta) \rightarrow 0$ as $m \rightarrow \infty$ for any fixed δ .

▷ In the non-iid setting: what if we had samples from bid vectors of players in past auctions, i.e.

if $F = F_1 \times \dots \times F_n$ and

$$S = \{ \vec{b}^1, \dots, \vec{b}^m \} \sim F^m$$

" $(b_{1,1}^1, \dots, b_{1,1}^m)$

Then can we construct a mechanism M_S s.t.

$$E_{S \sim F^m, \vec{v} \sim F} [Rev_{M_S}(\vec{v})] \geq OPT(F) - \epsilon(m)$$

or w.p. $1 - \delta$ over S :

$$E_{\vec{v} \sim F} [Rev_{M_S}(\vec{v})] \geq OPT(F) - \epsilon(m, \delta)$$

▷ This is exactly the type of question that PAC learning theory asks

Probably Approximately Correct (Valiant '84)

▷ The general setting of learning was first

analyzed by Vapnik and Chervonenkis' + 1
(Vapnik '82, '92, '95, '98)

- We have a hypothesis space H over which we want to learn and optimize (the space of dominant strategy truthful mechanisms)

- We have a distribution D over data points $z \in Z$
(the distribution F over bid/value vectors $\vec{b} \in [0, \#]^n$)

- We have an objective function in mind: (loss or reward) we will transition to losses
 $l: H \times Z \rightarrow \mathbb{R}$

(negative of the revenue of a mechanism M on a value vector \vec{b})

- We want to optimize expected loss

$$\inf_{h \in H} \mathbb{E}_{z \sim D} [l(h, z)] = \text{OPT}(D)$$

Known under many names: generalization error, risk, true error

- We have a set of iid samples from D , $S = \{z_1, \dots, z_n\}$
(the sample bid vectors).

Goal Given S find a hypothesis h_S

s.t. w.p. $1-\delta$ over the sample draws:

$$\underbrace{\mathbb{E}_{z \sim D} [\ell(h_S, z)]}_{\text{Denote as } L_D(h_S)} \leq \text{OPT}(D) + \epsilon$$

Sample Complexity: For any ϵ, δ is there an $m(\epsilon, \delta)$ s.t. \exists alg. that with $m > m(\epsilon, \delta)$ sample set S , we guarantee that w.p. $1-\delta$: $L_D(h_S) \leq \text{OPT}(D) + \epsilon$
 $m(\epsilon, \delta)$ is the sample complexity of the problem.

Ideally, we want it to be $\text{poly}(\frac{1}{\epsilon}, \log(1/\delta))$
Ideally, we want the algorithm to be poly time computable
→ Valiant's definition. 184

▷ Now we will just ask: is $m(\epsilon, \delta)$ finite for all ϵ, δ . Then we will say that the problem is learnable.

Candidate algorithm

- Since samples are drawn from the target distribution, the average loss on the samples is a good proxy for the expected loss:

$$L_S(h) = \frac{1}{n} \sum_{t=1}^n l(h, z_t)$$

- So why not just output the hypothesis h_S which minimizes $L_S(h)$

$$h_S = \operatorname{arg\,inf}_{h \in H} L_S(h)$$

- If $L_S(\cdot)$ is close to $L_D(\cdot)$
 - loss on samples
 - training error
 - empirical loss
- loss on true dist.
- generalization error
- true error

then $L_D(h_S)$ should also be small

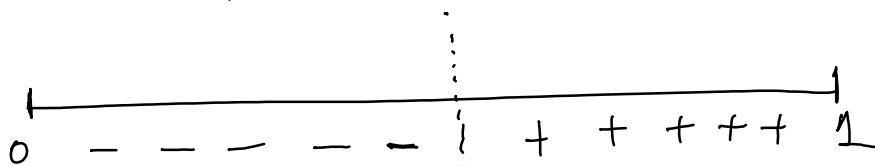
- This algorithm is called Empirical Risk Minimization or ERM.

- It is equivalent to the Follow-the-Leader Algorithm!

ERM gone bad

- Let's look at a classification problem in 1-d.

$$z = (x, y) \quad x \in [0, 1], \quad y \in \{0, 1\}$$



- $x \sim U(0, 1)$ $y|x$ is 1 if $x \geq \frac{1}{2}$ and 0 o.w.
- H : space of all mappings from $[0, 1] \rightarrow \{0, 1\}$
- loss: Prediction error; $l(h, (x, y)) = \mathbb{1}\{h(x) \neq y\}$

- ERM:

$$h_S = \arg \min_{h \in H} \frac{1}{m} \sum_{t=1}^m \mathbb{1}\{h(x_t) \neq y_t\}$$

- For instance: (memoization)
$$h_S(x) = \begin{cases} y_i & \text{if } \exists i \in [m] \text{ s.t. } x_i = x \\ 0 & \text{o.w.} \end{cases}$$

- Obviously: $L_S(h_S) = 0$ (zero loss on sample)

- But:

$$L_D(h_S) = \mathbb{E}_{(x,y) \sim D} [l(h_S, (x,y))] = \frac{1}{2}$$

since h_S predicts $\mathbb{1}$ only on a finite (zero-measure) x 's

- So we will never achieve $\epsilon < \frac{1}{2}$ error no matter how large m is!!

- ERM fails! **Over-fitting to samples!**

► Is it a problem of ERM??

- Not in this case.

- You can show that if the space X where the x "lives" is infinite, then the hypothesis space H of all functions from $X \rightarrow \{0,1\}$ is not learnable by any algorithm.

- We need to restrict H !!

- For instance: if H was the space of all threshold functions,

$$H = \left\{ f_\theta(x) = \begin{cases} 1 & \text{if } x \geq \theta \\ 0 & \text{o.w.} \end{cases} \text{ for } \theta \in [0,2] \right\}$$

Then this would fix our problem (see

at the end of next class).

Learnability from uniform convergence

▷ More generally the problem was that $|L_S(h) - L_D(h)|$ was not uniformly close for all h .

Uniform Convergence

- Suppose that $\forall \varepsilon, \delta : \exists m_U(\varepsilon, \delta)$ s.t.
 $\forall m > m_U(\varepsilon, \delta)$, w.p. $1 - \delta$

$$\sup_{h \in H} |L_S(h) - L_D(h)| \leq \varepsilon$$

- i.e. we know that the worst case error of the empirical loss goes to zero uniformly over all hypotheses.

Thm Uniform Convergence \Rightarrow Learnability

Pf

Let: $h_D^* = \arg \inf_{h \in H} L_D(h)$

Consider $m \geq m_U(\frac{\varepsilon}{2}, \delta)$, then w.p. $1 - \delta$

$$L_D(h_S) \leq L_S(h_S) + \frac{\varepsilon}{2} \leq L_S(h_D^*) + \frac{\varepsilon}{2}$$

$$L_D(h_S) \leq L_S(h_S) + \frac{\epsilon}{2} \leq L_S(h_D^*) \leq L_D(h_D^*) + \frac{\epsilon}{2} + \frac{\epsilon}{2}$$

↑
this
is where
our example
failed
Empirical loss deceptively
small.

Finite H and Uniform Convergence

Thm If H is finite and $l(h, \delta) \in [\alpha, b]$

$$m_U(\epsilon, \delta) \leq \left\lceil \frac{\log(2|H|/\delta)(b-\alpha)^2}{2\epsilon^2} \right\rceil$$

so H is learnable. and with sample complexity $\text{poly}(\frac{1}{\epsilon}, \log(1/\delta), \log(|H|))$

Pf] For a given m , we want to show:

$$P\left(\sup_{h \in H} |L_S(h) - L_D(h)| \leq \epsilon\right) \geq 1 - \delta$$

Equiv.

$$P\left(\sup_{h \in H} |L_S(h) - L_D(h)| > \epsilon\right) < \delta$$

Equiv

Equiv

$$P(\exists h \in \mathcal{H} : |L_S(h) - L_D(h)| > \epsilon) < \delta$$

By union bound

$$P(\exists h \in \mathcal{H} : |L_S(h) - L_D(h)| > \epsilon) \leq$$

$$\sum_{h \in \mathcal{H}} \underbrace{P(|L_S(h) - L_D(h)| > \epsilon)}_{\text{Let's bound this!}}$$

For any fixed h:

$$L_D(h) = \mathbb{E}_{z \sim D} [l(h, z)]$$

$$L_S(h) = \frac{1}{m} \sum_{t=1}^m \underbrace{l(h, z_t)}_{\text{Random variable with}}$$

$$\mathbb{E}_{z_t \sim D} [l(h, z_t)] = L_D(h).$$

• So $L_S(h)$ is the average of i.i.d. random variables each with expected value $L_D(h)$.

• We need to bound the prob. that this average deviates \times lot from its mean!

(similar to Law of Large Numbers but not asymptotic).

Hoeffding's Inequality

Let $\theta_1, \dots, \theta_m$ be i.i.d random variables with $\theta_i \in [a, b]$ and $E[\theta_i] = \mu$. Then

$$Pr \left[\left| \frac{1}{m} \sum_{t=1}^m \theta_t - \mu \right| > \varepsilon \right] \leq 2 \exp \left(- \frac{2m\varepsilon^2}{(b-a)^2} \right)$$

(one of the most useful inequalities in any large deviation analysis)

Back to main thm. we get:

$$P \left(\left| L_S(h) - L_D(h) \right| > \varepsilon \right) \leq 2 \exp \left(- \frac{2m\varepsilon^2}{(b-a)^2} \right)$$

Pickind $m = \frac{\log(2|H|/\delta) (b-a)^2}{2\varepsilon^2}$

above is at most: $\frac{2\varepsilon^2}{|H|}$

So when we sum over H we get at most δ . □

Bonus pf of Hoeffding

Bonus pf of Hoeffding

Consider $z_i = \sigma_i - \mu \in [\alpha, b]$ Show:

$$P\left(\left|\frac{1}{m} \sum_{+} z_{+}\right| > \varepsilon\right) \leq 2 \exp\left(-\frac{2m\varepsilon^2}{(b-\alpha)^2}\right)$$

$$\begin{aligned} P\left(\frac{1}{m} \sum_{+} z_{+} > \varepsilon\right) &= P\left(e^{\frac{\lambda}{m} \sum_{+} z_{+}} > e^{\lambda \varepsilon}\right) \\ &\leq \frac{\mathbb{E}\left[e^{\frac{\lambda}{m} \sum_{+} z_{+}}\right]}{e^{\lambda \varepsilon}} \quad (\text{Markov}) \\ &= \frac{\mathbb{E}\left[\prod e^{\frac{\lambda}{m} z_{+}}\right]}{e^{\lambda \varepsilon}} \\ &= \frac{\prod \mathbb{E}\left[e^{\frac{\lambda}{m} z_{+}}\right]}{e^{\lambda \varepsilon}} \quad (\text{independence}) \end{aligned}$$

For any $X \in [\alpha, b]$: $\mathbb{E}\left[e^X\right] \leq e^{\frac{(b-\alpha)^2}{8}}$

Apply to $X = \frac{\lambda}{m} z_{+}$:

$$\leq \frac{e^{\frac{\lambda^2 (b-\alpha)^2}{8m}}}{e^{\lambda \varepsilon}}$$

$$= e^{\frac{\lambda^2 (b-\alpha)^2}{8m} - \lambda \varepsilon}$$

Pick $\lambda = \frac{4m\varepsilon}{(b-\alpha)^2}$?

Pick $\eta = \frac{4m\varepsilon}{(b-a)^2}$
 $\leq e^{-\frac{2m\varepsilon^2}{(b-a)^2}}$



Intuition! e^X random variable
 punishes very large deviations of
 X above zero. So if we can
 upper bound $\mathbb{E}[e^X]$ we should
 be able to show great concentration
 bounds for X .

Independence decouples $\mathbb{E}[e^X] = (\mathbb{E}[e^{\frac{1}{m}X_i}])^m$
 if $X = \frac{1}{m} \sum_i X_i$
 So a bound on $\mathbb{E}[e^{X_i}]$ is amplified
 by raising it to the m -th power.