

## Lecture 24: Statistical Queries and Random Classification Noise

May 5th, 2020

Lecturer: Nika Haghtalab

Readings: N/A

Scribe: Aaron Tucker

**Last time**

- Defined Random Classification Noise (RCN)
- Defined statistical query model (SQ)

STAT( $h^*$ ,  $\mathcal{D}$ ) gets queries of the form  $(\psi, \tau)$ , where  $\psi : \mathcal{X} \times \{0, 1\} \rightarrow \{0, 1\}$  and  $0 < \tau < 1$ . The answer to a query is  $v$  such that

$$\left| v - \mathbf{E} \left[ \psi(x, h^*(x)) \right] \right| \leq \tau.$$

**Today**

Formally see the connection between the RCN and STAT model, and how learnability in these two settings are related.

**1 Warmup: Simulating statistical queries**

**Theorem 1.1.** *If  $\mathcal{H}$  is efficiently learnable in the statistical query model, it is also efficiently learnable in the realizable PAC model. In particular, if we use  $M$  queries each with a tolerance of at least  $\tau > 0$  to learn a hypothesis of  $\text{err}_{\mathcal{D}}(h) \leq \epsilon$ , then  $m$  samples are sufficient to get  $\text{err}_{\mathcal{D}}(h) \leq \epsilon$  with probability  $1 - \delta$  for*

$$m \in O \left( \frac{1}{\tau^2} \log \left( \frac{M}{\delta} \right) \right).$$

*Proof.* Consider using  $m$  samples  $\{(x_i, y_i)\}_{i=1}^m$  to learn a single query  $(\psi, \tau)$  using the empirical estimate

$$v = \frac{1}{m} \sum_{i=1}^m \mathbf{1}(\psi(x_i, y_i) = 1).$$

Note that  $\frac{1}{m} \sum \mathbf{1}(\psi(x_i, y_i) = 1)$  is the empirical estimate of  $\Pr(\psi(x, h^*(x)) = 1)$ , which is the same as  $\mathbf{E}(\psi(x, h^*(x)))$  because  $\psi$  has range  $\{0, 1\}$ . We can use the Hoeffding bound to show that

$$\left| \frac{1}{m} \sum_{i=1}^m \left( \mathbf{1}(\psi(x_i, y_i) = 1) \right) - \mathbf{E} \left[ \psi(x, h^*(x)) \right] \right| \leq \exp \left( \frac{-m\tau^2}{2} \right) \leq \delta.$$

If we have  $M$  queries, then we need to use a union bound to ensure that all of them are true with probability  $\delta$ , which means we need to solve

$$\exp\left(\frac{-m\tau^2}{2}\right) \leq \delta/M.$$

Solving, we see that if we can answer all  $M$  queries with  $\tau$  accuracy with probability  $1 - \delta$  as long as

$$m \geq \frac{1}{\tau^2} \log\left(\frac{M}{\delta}\right),$$

as desired. □

## 2 SQ-Learnability and RCN-Learnability

STAT-learnability implies learnability in the PAC with random classification noise model.

**Theorem 2.1.** *If  $\mathcal{H}$  is efficiently learnable in the statistical query model, it is also efficiently learnable in the PAC model with random classification noise. In particular, if we use  $M$  queries each with a tolerance of at least  $\tau > 0$  to learn a hypothesis of  $\text{err}_{\mathcal{D}}(h) \leq \epsilon$ , then  $m$  samples are sufficient to get  $\text{err}_{\mathcal{D}}(h) \leq \epsilon$  with probability  $1 - \delta$  for*

$$m \in O\left(\frac{1}{\tau^2(1-2\eta)^2} \ln\left(\frac{M}{\delta}\right)\right).$$

*Remark 2.2.* For ease of this next proof, instead of  $\psi : \mathcal{X} \times \{0, 1\} \rightarrow \{0, 1\}$  we use

$$\psi : \mathcal{X} \times \{-1, 1\} \rightarrow \{-1, 1\}.$$

This doesn't change the expressivity of the model, but makes the math easier for us.

### 2.1 Simplifying $\mathbf{E}_{x \sim \mathcal{D}} \left[ \psi(x, h^*(x)) \right]$

For any query  $\psi$ , we want to separate it into a part that cares only about the marginal distributions (but not the labels), and a part that also cares about the labels. This is similar to our previous proof for monotone ands – there were parts that cared about the relative frequency of 0 and 1 but not about the label or the noise, and other parts that cared about the labels given that certain features were set to 1.

$$\begin{aligned} \mathbf{E}_{x \sim \mathcal{D}} \left[ \psi(x, h^*(x)) \right] &= \mathbf{E}_{x \sim \mathcal{D}} \left[ \psi(x, 1) \cdot \mathbf{1}(h^*(x) = 1) \right] + \mathbf{E}_{x \sim \mathcal{D}} \left[ \psi(x, -1) \cdot \mathbf{1}(h^*(x) = -1) \right] \\ &= \mathbf{E}_{x \sim \mathcal{D}} \left[ \psi(x, 1) \cdot \frac{1 + h^*(x)}{2} \right] + \mathbf{E}_{x \sim \mathcal{D}} \left[ \psi(x, -1) \cdot \frac{1 - h^*(x)}{2} \right] \\ &= \underbrace{\frac{1}{2} \mathbf{E}_{x \sim \mathcal{D}} [\psi(x, 1) + \psi(x, -1)]}_{\text{target independent}} + \underbrace{\frac{1}{2} \mathbf{E}_{x \sim \mathcal{D}} [\psi(x, 1) h^*(x) - \psi(x, -1) h^*(x)]}_{\text{correlation of } h^*(x) \text{ with some other function}} \end{aligned}$$

This works because

$$\mathbf{1}(h^*(x) = 1) = \frac{1 + h^*(x)}{2} \text{ and } \mathbf{1}(h^*(x) = -1) = \frac{1 - h^*(x)}{2}$$

This substitution lets us split any query  $\psi$  into a target independent term that cares only about the marginal distribution of  $x_i$ , and a correlational term that cares about the correlation between  $h^*(x)$  and some other function.

## 2.2 Target Independent and Correlational Queries

We can use these two types of queries without loss of generality.

**Definition 2.3** (Correlational queries). A query of the form  $\varphi : \mathcal{X} \rightarrow \{-1, 1\}$  and tolerance  $\tau \in (0, 1)$  which returns a  $v$  as follows is called a correlational query.

$$\left| \mathbf{E} [\varphi(x) \cdot h^*(x)] - v \right| \leq \tau$$

This gives us information about  $h^*$ .

**Definition 2.4** (Target independent queries). A query of the form  $\phi : \mathcal{X} \rightarrow \{-1, 1\}$  and tolerance  $\tau \in (0, 1)$  which returns a  $v$  as follows is called a target independent query.

$$\left| \mathbf{E} [\phi(x)] - v \right| \leq \tau$$

This is sufficient to answer questions about the marginal distribution.

## 2.3 Simulating Queries in the RCN Model

**Simple case: Target independent queries**

Target independent  $\phi : \mathcal{X} \rightarrow \{-1, 1\}$ . The noise has no effect, so given  $\{(x_1, y_1), \dots, (x_m, y_m)\}$ , just throw out the  $y_i$  and compute

$$v = \frac{1}{m} \sum_{i=1}^m \phi(x_i).$$

This is just the same as the warm up exercise, since throwing out  $y_i$  throws out the source of noise, and it is the same as in the case for deterministic labels.

$$\left| v - \mathbf{E} [\phi(x)] \right| \leq \tau.$$

Given  $\frac{1}{\tau^2} \log \left( \frac{M}{\delta} \right)$  data points, we succeed with probability  $1 - \delta$ .

## Correlational queries

Recall that  $\mathcal{D}(\eta)$  is the noisy version of  $\mathcal{D}$ , with random classification noise  $\eta$ . Answering the query  $\varphi$  means approximating  $\mathbf{E}_{x \sim \mathcal{D}} [\varphi(x) \cdot h^*(x)]$ . However we can only have direct access to  $\mathbf{E}_{(x,y) \sim \mathcal{D}(\eta)} [\varphi(x) \cdot y]$ . How can we relate these?

**Claim 2.5.**  $\mathbf{E}_{x \sim \mathcal{D}} [\varphi(x) \cdot h^*(x)]$  and  $\mathbf{E}_{(x,y) \sim \mathcal{D}(\eta)} [\varphi(x) \cdot y]$  are related by:

$$\mathbf{E}_{(x,y) \sim \mathcal{D}(\eta)} [\varphi(x) \cdot y] = (1 - 2\eta) \mathbf{E}_{x \sim \mathcal{D}} [\varphi(x) \cdot h^*(x)].$$

*Proof.* Because  $\varphi$  has range  $\{-1, 1\}$ , we can say that

$$\mathbf{E}_{(x,y) \sim \mathcal{D}(\eta)} [\varphi(x) \cdot y] = \Pr_{\mathcal{D}(\eta)} [\varphi(x) = y] - \Pr_{\mathcal{D}(\eta)} [\varphi(x) \neq y]$$

The first term is the probability that  $\varphi(x)$  agrees with  $y$ . That can either happen because  $\varphi(x)$  agrees with  $h^*(x)$  and  $y$  wasn't flipped, or because  $\varphi(x)$  disagrees with  $h^*(x)$ , but  $y$  was flipped.

$$\begin{aligned} \Pr_{\mathcal{D}(\eta)} [\varphi(x) = y] &= \Pr(\varphi(x) = h^*(x)) \cdot \Pr(h^*(x) = y) + \Pr(\varphi(x) \neq h^*(x)) \cdot \Pr(h^*(x) \neq y) \\ &= \Pr(\varphi(x) = h^*(x)) (1 - \eta) + \left(1 - \Pr(\varphi(x) = h^*(x))\right) \eta \\ &= (1 - 2\eta) \Pr(\varphi(x) = h^*(x)) + \eta. \end{aligned}$$

The second term is the probability that  $\varphi(x)$  disagrees with  $y$ . That can either happen because  $\varphi(x)$  disagrees with  $h^*(x)$  and  $y$  wasn't flipped, or because  $\varphi(x)$  agrees with  $h^*(x)$  but  $y$  was flipped.

$$\begin{aligned} \Pr_{\mathcal{D}(\eta)} [\varphi(x) \neq y] &= \Pr(\varphi(x) \neq h^*(x)) \cdot \Pr(h^*(x) = y) + \Pr(\varphi(x) = h^*(x)) \cdot \Pr(h^*(x) \neq y) \\ &= \Pr(\varphi(x) \neq h^*(x)) (1 - \eta) + \left(1 - \Pr(\varphi(x) \neq h^*(x))\right) \eta \\ &= (1 - 2\eta) \Pr(\varphi(x) = h^*(x)) + \eta. \end{aligned}$$

Notice that both terms are now expressed in terms of the probability of how  $\varphi(x)$  relates to the true  $h^*(x)$ , rather than to the noisy  $y$ .

$$\begin{aligned} \mathbf{E}_{(x,y) \sim \mathcal{D}(\eta)} [\varphi(x) \cdot y] &= \Pr_{\mathcal{D}(\eta)} [\varphi(x) = y] - \Pr_{\mathcal{D}(\eta)} [\varphi(x) \neq y] \\ &= (1 - 2\eta) \Pr(\varphi(x) = h^*(x)) + \eta - \left((1 - 2\eta) \Pr(\varphi(x) = h^*(x)) + \eta\right) \\ &= (1 - 2\eta) \left(\Pr(\varphi(x) = h^*(x)) - \Pr(\varphi(x) = h^*(x))\right) \\ &= (1 - 2\eta) \mathbf{E}_{x \sim \mathcal{D}} [\varphi(x) \cdot h^*(x)] \end{aligned}$$

as desired. □

### Number of samples needed to estimate $v$

If we want to estimate a correlational query  $v = \mathbf{E}_{\mathcal{D}} [\varphi(x) \cdot h^*(x)]$  to accuracy  $\tau$ , we should estimate  $v' = \mathbf{E}_{\mathcal{D}(\eta)} [\varphi(x) \cdot h^*(x)]$  to accuracy  $\tau'$ , where  $\tau'$  is a parameter to be set later. Consider  $m$  samples  $\{(x_i, y_i)\}_{i=1}^m$  from the RCN oracle  $\mathbb{E}X^\eta(h^*, \mathcal{D})$ , and compute  $v'$  as

$$v' = \frac{1}{m} \sum_{i=1}^m \varphi(x_i) \cdot y_i.$$

setting  $m$  such that

$$\left| v' - \mathbf{E}_{\mathcal{D}(\eta)} [\varphi(x) \cdot y] \right| \leq \tau'.$$

Setting  $v = \frac{v'}{1-2\eta}$  as our estimate for  $\mathbf{E}_{\mathcal{D}} [\varphi(x) \cdot h^*(x)]$ , we get

$$\left| v - \mathbf{E}_{\mathcal{D}} [\varphi(x) \cdot h^*(x)] \right| = \left| \frac{v'}{1-2\eta} - \mathbf{E}_{\mathcal{D}} [\varphi(x) \cdot h^*(x)] \right| \tag{1}$$

$$= \left| \frac{v'}{1-2\eta} - \frac{1}{1-2\eta} \mathbf{E}_{\mathcal{D}(\eta)} [\varphi(x) \cdot y] + \frac{1}{1-2\eta} \mathbf{E}_{\mathcal{D}(\eta)} [\varphi(x) \cdot y] - \mathbf{E}_{\mathcal{D}} [\varphi(x) \cdot h^*(x)] \right| \tag{2}$$

$$= \frac{1}{1-2\eta} \left| v' - \mathbf{E}_{\mathcal{D}(\eta)} [\varphi(x) \cdot y] \right| \tag{3}$$

$$\leq \frac{1}{1-2\eta} \tau' \tag{4}$$

$$\tag{5}$$

At step 2 we add and subtract  $\frac{1}{1-2\eta} \mathbf{E}_{\mathcal{D}(\eta)} [\varphi(x) \cdot y]$ . We can cancel out the last two terms of line (2) because of [Claim 2.5](#). Setting  $\tau' = (1-2\eta)\tau$ , we get

$$\left| v - \mathbf{E}_{\mathcal{D}} [\varphi(x) \cdot h^*(x)] \right| \leq \tau$$

We know from [Theorem 1.1](#) that the sample complexity of learning  $v'$  is  $\frac{1}{\tau'^2} \log\left(\frac{M}{\delta}\right)$ , so if we substitute in our expression for  $\tau'$  we get a sample complexity of

$$\frac{1}{\tau^2(1-\eta)^2} \log\left(\frac{M}{\delta}\right).$$

### Conclusion

This shows that anything you can learn efficiently in the statistical query model, you can learn in the PAC+RCN model, and if you needed  $M$  queries of tolerance  $\tau$ , then all you need is  $m \in O\left(\frac{1}{\tau^2(1-\eta)^2} \log\left(\frac{M}{\delta}\right)\right)$  data points to learn in the PAC+RCN model.

### 3 Discussion

#### What if $\eta$ is unknown?

We can guess and try it. For example, if we have an effective upper bound  $\eta' \geq \eta$ , and  $\eta' - \eta = \Delta$  for a small  $\Delta$ , then as long as  $\Delta$  is small compared to  $\tau$  it won't really affect our accuracy. We just need a few more samples.

So, we can try our algorithm with  $\eta' \in \{0, \Delta, 2\Delta, \dots, 1/2\}$ . We know that one of these  $\eta'$  is close enough to be okay, and we can look at the error that we observe.

#### Why are correlational queries interesting?

- They give us insight into how much of the difficulty of the learning task is coming from its labels, as opposed to coming from its marginal distribution. If I knew  $\mathcal{D}$  the marginal, the only type of queries that you need to ask are correlational. This gives rise to a statistical query dimension.
- $\text{SQ-Dim}_{\mathcal{D}}(\mathcal{H}) = d$  if I have functions  $f_1, f_2, \dots, f_d \in \mathcal{H}$ , such that for all  $i \neq j$  we have  $\mathbf{E}_{\mathcal{D}} [f_i(x) \cdot f_j(x)] \leq \frac{1}{d} \approx 0$ . This lets us say that  $f_i$  are effectively an orthogonal basis vectors for  $\mathcal{H}$ . If you have  $\text{SQ-Dim}_{\mathcal{D}}(\mathcal{H}) = d$ , then you need  $\Omega(\tau^2 d)$  queries of  $\tau$  accuracy. Since the  $f_i$  are not correlated with each other, any query  $\varphi$  will be measuring one of them, but not the rest of them. So there has to be many of these queries in order to learn. You will formalize this in HW5.