

## Lecture 9

Instructor: *Jess Sorrell*Scribe: *Jess Sorrell*

**Acknowledgements.** Much of this material (and the material for the next few weeks) is lifted wholesale from the course notes of Aaron Roth and Adam Smith, available at

<https://www.adaptivedataanalysis.com>

Domain  $\mathcal{X} = \{0, 1\}^d$ ,  $\mathcal{Y} = \{0, 1\}$ .

## Overfitting with “natural” adaptive SQs

---

**Algorithm 1** Query learner

Inputs/Parameters: Sample  $S \sim D^m$

---

```

1:  $P = \emptyset$ 
2: for  $i \in [d]$  do
3:    $\phi_i(x, y) = \begin{cases} 1, & x_i = y \\ 0, & \text{o.w.} \end{cases}$ 
4:    $a_i \leftarrow \frac{1}{m} \sum_{(x,y) \in S} [\phi(x, y)]$ 
5:   if  $a_i \geq \frac{1}{2} + \frac{1}{\sqrt{m}}$  then
6:      $P = P \cup i$ 
7:   end if
8: return  $f(x) = \lfloor \frac{1}{|P|} \sum_{i \in P} x_i \rfloor$ 
9: end for

```

---

**Claim 0.1.** When  $D$  is the uniform distribution over  $\mathcal{X} \times \mathcal{Y}$ ,  $\exists$  constant  $c$  such that with probability at least  $1 - \delta$ , if  $d \geq c \max\{m, \log(1/\delta)\}$ :

$$|acc_S(f) - acc_D(f)| \geq .49$$

Compare to the accuracy guarantee we have for non-adaptive statistical queries, from which we would expect

$$|acc_S(f) - acc_D(f)| \in O\left(\sqrt{\frac{\log(d/\delta)}{m}}\right).$$

*Proof.* Let

$$X_i = \begin{cases} 1, & i \in P \\ 0, & \text{o.w.} \end{cases}$$

Then

$$\Pr_{S \sim D^m}[X_i = 1] = \Pr_{S \sim D^m} \left[ \frac{1}{m} \sum_{(x,y) \in S} \mathbb{1}[x_i = y] \geq \frac{1}{2} + \frac{1}{\sqrt{m}} \right]$$

Let  $A_i = \frac{1}{m} \sum_{(x,y) \in S} \mathbb{1}[x_i = y]$ .  $A_i$  is a binomial random variable with  $\mathbb{E}_D[A_i] = \frac{1}{2}$  and standard deviation  $\frac{1}{2\sqrt{m}}$ , and therefore

$$\Pr_D[X_i = 1] = \Pr_D[A_i \geq \frac{1}{2} + \frac{1}{\sqrt{m}}] \in \Omega(1).$$

Therefore, each  $i$  gets added to  $P$  with constant probability. It follows that

$$\mathbb{E}_{S \sim D^m}[|P|] = \sum_{i=1}^d \Pr[X_i = 1] = \Omega(d).$$

Recall what the Chernoff-Hoeffding inequality gives us for a sum of bounded r.v.'s  $X_i \in [a_i, b_i]$ , letting  $S_d = \sum_{i=1}^d X_i$ :

$$\Pr_{X_1, \dots, X_d}[S_d \leq \mathbb{E}[S_d] - t] \leq e^{-\frac{2t^2}{\sum_{i=1}^d (b_i - a_i)^2}}$$

applied to  $|P|$ , there is a  $t \in \Omega(d)$  such that we have

$$\begin{aligned} \Pr_{S \sim D^m}[|P| \in o(d)] &\leq \Pr_{X_1, \dots, X_d}[|P| \leq \mathbb{E}[|P|] - t] \\ &\leq e^{-\frac{2t^2}{\sum_{i=1}^d (b_i - a_i)^2}} \\ &= e^{-\frac{2t^2}{d}} \\ &\in e^{-\Omega(d)} \end{aligned}$$

So there exists some constant  $c_1$  such that so long as  $d > c_1 \log(1/\delta)$ ,

$$\Pr_{S \sim D^m}[|P| \in \Omega(d)] > 1 - \delta$$

Now let's see how this causes us to get an unreliable empirical estimate of  $acc_S(f)$  if we reuse the sample  $S$ . Let  $(x, y) \sim S$  be chosen uniformly at random.  $f(x) = y$  iff  $\sum_{i \in P} \mathbb{1}[x_i = y] \geq \frac{|P|}{2}$ . We have that for each  $i \in P$ ,  $\Pr[x_i = y] \geq \frac{1}{2} + \frac{1}{\sqrt{m}}$ , and so

$$\mathbb{E}_{(x,y) \sim S} \left[ \sum_{i \in P} \mathbb{1}[x_i = y] \right] \geq \frac{|P|}{2} + \frac{|P|}{\sqrt{m}}$$

Therefore  $f(x) = y$  unless  $\sum_{i \in P} \mathbb{1}[x_i = y]$  is less than its expectation by at least  $\frac{|P|}{\sqrt{m}}$ . Applying Chernoff-Hoeffding to the sum of random variables  $C = \sum_{i \in P} \mathbb{1}[x_i = y]$ , we have

$$\begin{aligned} 1 - acc_S(f) &= \Pr_{(x,y) \sim S}[f(x) \neq y] \\ &= \Pr_{(x,y) \sim S}[C \leq \mathbb{E}[C] - \frac{|P|}{\sqrt{m}}] \\ &\leq e^{-\frac{2|P|^2}{m|P|}} \\ &= e^{-\frac{2|P|}{m}} \end{aligned}$$

So if  $|P| > \frac{\ln(100)m}{2}$ ,  $acc_S(f) = .99$ . However,  $acc_D(f) = 1/2$ . We already showed that there exists  $c_1$  such that  $|P| \in \Omega(d)$  except with probability  $\delta$ , so long as  $d > c_1 \log(1/\delta)$ . Therefore there exists  $c_2$  such that so long as  $d > c_2 m$ ,  $|P| > \frac{\ln(100)m}{2}$ , and  $acc_S(f) = .99$ . Letting  $c = \max\{c_1, c_2\}$ , it follows that there exists a  $c$  such that with probability at least  $1 - \delta$ , if  $d \geq c \max\{m, \log(1/\delta)\}$ :

$$|acc_S(f) - acc_D(f)| \geq .49$$

□

## Observations

- The argument above still goes through when we don't answer queries with an exact empirical estimate, but instead add some noise on the order  $o(\frac{1}{\sqrt{n}})$  to the estimate.
- We could have done a similar analysis using only the first  $k - 1$  of  $d$  features. Redoing the argument using only  $k$  statistical queries instead of  $d + 1$  gives a bound of

$$|acc_S(f) - acc_D(f)| \in \Omega(\sqrt{km}).$$

So the best confidence interval we can hope for with  $k$  adaptive statistical queries, answered by empirical estimate over reused data, is  $O(\sqrt{\frac{k}{m}})$ . Recall that for  $k$  non-adaptive statistical queries, our bound was  $O(\sqrt{\frac{\log k}{m}})$ .

- If we want a confidence interval of  $\varepsilon$ , reusing data in this way doesn't save us anything, since we would need  $m \in \Omega(\frac{k}{\varepsilon^2})$  samples... which is  $k$  times what we would need for a single statistical query.
- Could we have made our adaptive algorithm non-adaptive? The first  $d$  queries will non-adaptive, so what if we just committed to estimating the error of every possible  $f$  we *might* have constructed in our algorithm, rather than the one that was chosen *after* looking at the results of our  $d$  non-adaptive queries. Would this give a better bound? There are  $2^d$  many subsets of  $d$  variables that could be included in  $P$ , and therefore  $2^d$  different functions  $f$ . So we would need to make  $O(2^d)$  statistical queries, giving a bound of  $O(\sqrt{\frac{\log 2^d}{m}}) = O(\sqrt{\frac{d}{m}})$ , so no better than the adaptive version.
- Do replicable SQs help? Since the first  $d$  queries are non-adaptive, we know that so long as we use a large enough sample, we can guarantee

$$\Pr_{S_1, S_2, r} [f_{S_1}^r = f_{S_2}^r] > 1 - \rho$$

where  $f_{S_i}^r = A(S_i; r)$ . It follows that

$$\begin{aligned}
\Pr_{S_1, r} [acc_{S_1}(f_{S_1}^r) \geq \frac{1}{2} + \tau] &= \Pr_{S_1, S_2, r} [acc_{S_1}(f_{S_2}^r) \geq \frac{1}{2} + \tau \mid f_{S_2}^r = f_{S_1}^r] \cdot \Pr_{S_1, S_2, r} [f_{S_2}^r = f_{S_1}^r] \\
&\quad + \Pr_{S_1, S_2, r} [acc_{S_1}(f_{S_1}^r) \geq \frac{1}{2} + \tau \mid f_{S_2}^r = f_{S_1}^r] \cdot \Pr_{S_1, S_2, r} [f_{S_2}^r \neq f_{S_1}^r] \\
&\leq \Pr_{S_1, S_2, r} [acc_{S_1}(f_{S_2}^r) \geq \frac{1}{2} + \tau \mid f_{S_2}^r = f_{S_1}^r] \cdot \Pr_{S_1, S_2, r} [f_{S_2}^r = f_{S_1}^r] + \rho \\
&\leq \Pr_{S_1, S_2, r} [acc_{S_1}(f_{S_2}^r) \geq \frac{1}{2} + \tau] + \rho \\
&= \Pr_{S_1, S_2, r} [acc_{S_1}(f_{S_2}^r) - \mathbb{E}_{S_1, S_2, r} [acc_{S_1}(f_{S_2}^r)] \geq \tau] + \rho \\
&\leq e^{-2\tau^2 m} + \rho \\
&\in O(\rho)
\end{aligned}$$

so long as we take  $m \in \Omega(\frac{\log 1/\rho}{\tau^2})$ .

However, to ensure that  $\Pr_{S_1, S_2, r} [f_{S_2}^r \neq f_{S_1}^r] \leq \rho$ , we need to make  $d$  non-adaptive replicable statistical queries with  $\rho' = \rho/d$ , so we need  $O(\frac{d^2}{\tau^2 \rho^2})$  samples. Which is already worse than resampling!