## Lecture 10: PAC learning from Fourier concentration

*Lecturer: Eshan Chattopadhyay*        *Scribe: Kevin Hua*

# 1 Low degree algorithm

Recall that there are two models of *PAC learning* that we are considering in this class.

And those two are $\begin{cases} \text{Random Example Model: we are given} \quad \{(x_i, f(x_i))\}_{i=1}^t, x \sim \{-1, 1\}^n \\ \text{Query Model: we are given an oracle that answers any query of } f \end{cases}$

Our goal is to learn any $f \in \mathcal{F}$, where $\mathcal{F}$ is a *concept class*.

In other words, in any of these two models, we want to design a learning algorithm $\mathcal{A}$ that will output $h : \{-1, 1\}^n \to \{-1, 1\}$, such that, with high probability, $\text{dist}(h, f) = \Pr[h \neq f] \leq \epsilon$.

**Theorem 1.1.** *Suppose for any $f \in \mathcal{F}$, $\boldsymbol{W}^{\geq k}[f] \leq \underset{\underset{|s| \geq k}{\overset{\shortparallel}{\sum} \hat{f}(s)^2}}{\eta} \underset{\underset{assume\ \eta \leq \epsilon/2}{\uparrow}}{}$ . Then, $\mathcal{F}$ is PAC-learnable with*

$poly(n^k, 1/\epsilon)$ *samples.*

The first step of proving this theorem is to approximate the low-degree fourier coefficients. To do this, we need the subroutine below.

Subroutine: Estimate Fourier coefficients using $FOURIER(s)_{f, \epsilon_1, \delta_1}$

---
**Algorithm 1:** $FOURIER(S)_{f, \epsilon_1, \delta_1}$

---
$$\overset{\text{as defined in last lecture}}{\downarrow}$$

**1** Output $\quad \widetilde{\hat{f}(S)} \quad \in \mathbb{R}$, such that, with probability at least $1 - \delta_1$,

$|\widetilde{\hat{f}(S)} - \hat{f}(S)| < \epsilon_1$.

---

Note that we have independent samples $\{(x_i, f(x_i))\}_{i=1}^t$. And we supply $FOURIER(s)_{f, \epsilon_1, \delta_1}$ with these samples to generate $\widetilde{\hat{f}(S)}$.

With a handy algorithm estimating any Fourier coefficients, our final algorithm is as follows. The intuition is simple: we don't have the leisure to estimate every Fourier coefficient as there are exponentially many of them. However, given we know the high-degree Fourier coefficients are tiny, we might well get by without estimating them.

---
**Algorithm 2:** $\mathcal{A}$[low degree algorithm]

---
**1** Estimate $\hat{f}(S), \forall S \subseteq [n], |S| \leq k$ using $FOURIER(S)_{f, \epsilon_1, \delta_1}$

**2** output $\text{sign}(h(x))$, where $\quad h(x) = \sum_{S \subseteq [n], |S| \leq k} \widetilde{\hat{f}(S)} \chi_S(x).$

---

Analysis:

Assume all estimates in step 1 are "good" $--|\widetilde{\hat{f}(S)} - \hat{f}(S)| \le \epsilon_1, \forall S \subseteq [n]$

Note that by Union Bound, $\Pr[\exists \text{a bad estimate}] \le \delta_1\binom{n}{\le k}$. At the end, we are going to pick $\delta_1$ so that this probability is some constant strictly less than 1.

Goal: we want to show $\text{dist}(\text{sign}(h), f) = \Pr[\text{sign}(h(x)) \ne f(x)]$ is small.

Let $g(x) = f(x) - h(x)$. Note that when $\text{sign}(h(x)) \ne f(x)$, $|g(x)| = |f(x) - h(x)| > 1$

Thus, we have that $\text{dist}(\text{sign}(h), f) = \Pr[\text{sign}(h(x)) \ne f(x)] \le ||g||_2 = \mathbf{E}[(f - h)^2]$.

By Parseval's Theorem,

$$||g||_2^2 = \sum_{s \subseteq [n]} (\hat{g}(S))^2$$

$$= \sum_{s:|S| \le k} \hat{g}(S)^2 + \mathbf{W}^{\ge k+1}[g]$$

As $g(x) = f(x) - h(x)$, $\hat{g}(S) = \hat{f}(S) - \hat{h}(S)$.

For $|S| \le k$, by definition,

$$|\hat{g}(S)| = |\hat{f}(S) - \hat{h}(S)|$$

$$= |\hat{f}(S) - \widetilde{\hat{f}(S)}|$$

$$\le \epsilon_1 \quad \text{(by assumption that the estimate } \widetilde{\hat{f}(S)} \text{ is "good")}$$

For $\mathbf{W}^{\ge k+1}[g]$, as $h(S) = 0$ for $|S| \ge k+1$, we have $\mathbf{W}^{k+1}[g] = \mathbf{W}^{k+1}[f] \le \eta \le \epsilon/2$.

Summing up the upper bounds for $\sum_{|S| \le k} \hat{g}(S)^2$ and the one for $\mathbf{W}^{k+1}[g]$, we get

$$\text{dist}(\text{sign}(h), f) \le ||g||_2^2 \le \binom{n}{\le k} \epsilon_1^2 + \epsilon/2$$

$$\le \binom{n}{\le k} \epsilon_1 + \epsilon/2 \quad \text{(for } \epsilon_1 \le 1\text{) with probability} \ge 1 - \delta_1 \binom{n}{\le k}.$$

Now, setting $\epsilon_1 = \frac{\epsilon}{2\binom{n}{\le k}} \sim O(\frac{\epsilon}{n^k})$, $\delta_1 = O(\frac{1}{n^k})$, we get that with some positive constant probability, $\mathcal{A}$ outputs a $h$ such that $\text{dist}(h, f) \le \epsilon$. $\qquad\square$

Note that sampling complexity of FOURIER is $O(\frac{kn^{2k}}{\epsilon^2} \log(n))$ from last time.

As we repeat FOURIER $O(\binom{n}{\le k})$ times in $\mathcal{A}$, we get the sampling complexity of $\mathcal{A}$ is $O(\frac{k \cdot n^{3k} \log n}{\epsilon^2})$.

## 2 Kushilevitz-Mansour Algorithm

**Theorem 2.1.** *Suppose $\mathcal{F}$ is $\eta$-concentrated as follows: for any $f \in \mathcal{F}$, $\exists L_f = \{s_1, \cdots, s_M\}$, s.t. $\sum_{s \in L_f} \hat{f}(s)^2 \ge 1 - \eta$. Then, $\mathcal{F}$ is $\epsilon$-PAC learnable in query model with sample complexity* $\text{poly}(n, \frac{1}{\epsilon}, M)$*, with $\eta \le \frac{\epsilon}{2}$.*

First, suppose $L_f$ is given to $\mathcal{A}$.

Then we can estimate $\hat{f}(S), \forall S \in L_f$ by outputting $\text{sign}(\sum_{s \in L_f} \widetilde{\hat{f}(s)} \chi_s(x))$.

Repeating the same analysis as for the low-degree algorithm, we can easily see that with $\text{poly}(n, \frac{1}{\epsilon}, M)$ samples, $\text{dist}(\text{sign}(\sum_{s \in L_f} \widetilde{\hat{f}(s)} \chi_s(x)), f) \leq \epsilon$. Therefore, we are done if we can find $L_f$(efficiently). GOOD NEWS: *Goldreich-Levin* algorithm find the $L_p$ for us!

## 3    Goldreich-Levin

**Theorem 3.1.** *With query access to $f$ and parameter $\delta$, Goldreich-Levin outputs a set $\tilde{L}_f$. With high probability, $\tilde{L}_f$ satisfies the following two constraints:*

*(1) if $S \in \tilde{L}_f$, $|\hat{f}(S)| \geq \frac{\delta}{2}$.*

*(2) if $|\hat{f}(S)| \geq \delta$, $S \in \tilde{L}_f$.*

<u>Notation</u>: For any $0 \leq k \leq n$, $S \subseteq [k]$

$$B_{k,S} := \{T \subseteq [n] : T \cap [k] = S\}$$

$$W(B_{k,S}) = \sum_{T \in B_{k,S}} \hat{f}(T)^2$$

Note that $|B_{k,S}| = 2^{n-k}$ and $B_{0,\emptyset} = 2^{[n]}$.

---

**Algorithm 3:** Goldreich-Levin Algorithm

---

**1** Set $\mathcal{B} = \{B_{0,\emptyset}\}$

    ↑ collection of buckets

**2 while** $\exists B \in \mathcal{B}$ *such that $B$ contains at least 2 sets* **do**

**3**     Suppose $B = B_{k,S}$. Remove $B$ from $\mathcal{B}$.

**4**     Let $B_0 = B_{k+1,S}$, $B_1 = B_{k+1,S \cup \{i+1\}}$.

**5**     Estimate $W(B_0)$ and $W(B_1)$ to accuracy $\frac{\delta^2}{4}$.

**6**     Add $B_i$ to $\mathcal{B}$ if $B_i$ is estimated with weight $\geq \frac{\delta^2}{2}$ (We call $B_i$ "heavy" in this case).
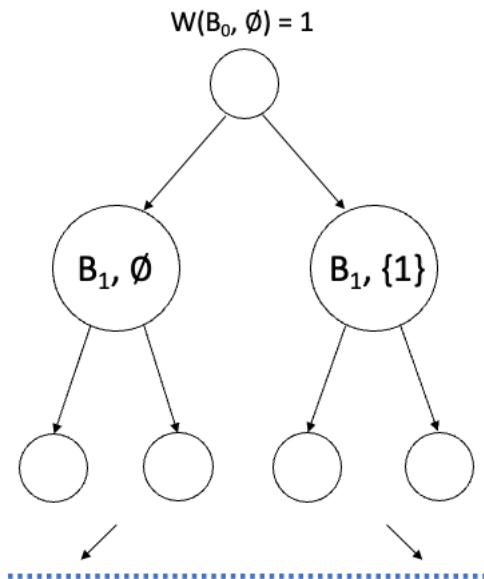
**7 end**

**8** Output $\mathcal{B}$.

---

Figure 1: representation of buckets as a complete binary tree

Question: How many heavy leaves are there in the trees that can be included as part of the output for *Goldreich-Levin*, as shown in Figure 1?

3Answer: $\frac{4}{\delta^2}$. Note that the leaves are disjoint buckets. Therefore, given the accuracy of the algorithm is $\frac{\delta^2}{4}$, there are at most $\frac{4}{\delta^2}$ leaves buckets that can be included in the output.