

Lecture 21

Lecturer: Madhu Sudan

Scribe: Paul Valiant

1 Trevisan's Extractor

The following is Trevisan's construction of an extractor. Consider a code $C : \{0, 1\}^n \rightarrow \{0, 1\}^N$ that is a $(\frac{1}{2} - \delta, L)$ soft-decision list decoder for some list size $L = \text{poly}(\frac{n}{\delta})$. Note that we have seen explicit constructions for codes with $N = \frac{n}{\delta^8}$ and $L \leq N$.

Recall that a (m, t, l, a) -design is a collection of sets S_1, \dots, S_m with $S_i \subset [t]$, $|S_i| = l$, and the condition that for every i, j , we have

$$|S_i \cap S_j| \leq a.$$

We note that Trevisan shows these designs exist if

$$t \geq \frac{l^2}{a} \exp\left(\frac{\log m}{a}\right).$$

For $l = \log_2 N$ we define an extractor $E : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$ as

$$E(x, y) = (C(x)_{y_1}, C(x)_{y_2}, \dots, C(x)_{y_m})$$

where y_i is the restriction of y to the set S_i , which is a string in $\{0, 1\}^l$ and can be interpreted to be an integer in $[N]$.

We provide a proof sketch for the following claim.

Claim 1 E is a (k, ϵ) extractor provided $\delta < \frac{\epsilon}{2m}$ and $k > m2^k + O(\log \frac{n}{\epsilon})$.

Proof (Sketch) Assume for the sake of contradiction that A ϵ -distinguishes a random $x \in X$ for some X with $|X| = 2^k$, namely

$$\Pr_{x \in X, y \in U_t} [A(E(x, y), y) = 1] > \Pr_{z \in U_m, y \in U_t} [A(z, y) = 1] + \epsilon,$$

where U_k represents the uniform distribution on strings in $\{0, 1\}^k$.

Using standard Markov arguments, one can show that there exists some set $X_{\frac{\epsilon}{2}} \subset X$ of size at least $\frac{\epsilon}{2}|X| = \frac{\epsilon}{2}2^k$ such that A $\frac{\epsilon}{2}$ -distinguishes every $x \in X_{\frac{\epsilon}{2}}$.

We now show that if one can $\frac{\epsilon}{2}$ -distinguish m bits of $E(x, y)$ from the uniform distribution, then one can $\frac{\epsilon}{2m}$ -distinguish a single bit from uniform. We use a technique called *hybridization*.

Define distributions D_0, \dots, D_m as follows: for fixed x and random $y \in U_t$ consider the m pseudo-random bits of $E(x, y) = (C(x)_{y_1}, \dots, C(x)_{y_m})$. Further, let $b_1 \dots b_m$ be m truly random bits. Define the distribution D_i by the distribution of

$$(C(x)_{y_1}, \dots, C(x)_{y_{i-1}}, b_{i+1}, \dots, b_m).$$

Let

$$p_i \stackrel{\text{def}}{=} \Pr_{(z, y) \in D_i} [A(z, y) = 1].$$

Note that by definition

$$p_m - p_0 > \frac{\epsilon}{2}.$$

Further, we may expand this into a telescoping sum as

$$\frac{\epsilon}{2} < p_m - p_0 = \sum_i p_i - p_{i-1},$$

which implies that for some i ,

$$p_{i+1} - p_i > \frac{\epsilon}{2m}.$$

Thus for each $x \in X_{\frac{\epsilon}{2}}$,

$$\mathbb{E}_{y,b} [A(C(x)_{y_1}, \dots, C(x)_{y_i}, b_{i+1}, \dots, b_m) - A(C(x)_{y_1}, \dots, C(x)_{y_{i-1}}, b_i, \dots, b_m)] > \frac{\epsilon}{2m},$$

from which we conclude that for any x there exists a suffix b_{i+1}, \dots, b_m such that the above equation holds. Further, for some specific suffix b_{i+1}, \dots, b_m the set of strings x for which the above equation holds $X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m}$, is at least $2^{-m} |X_{\frac{\epsilon}{2}}| \geq \frac{\epsilon}{2} 2^{k-m}$.

Similarly, the above bound is an average over all y , but we are concerned only with the digits of y that help determine the bit $C(x)_{y_i}$, namely those digits of y indexed by elements of the set S_i . Thus, as above, we conclude that for any x there exists an assignment to $y|_{[t]-S_i}$ such that we can distinguish x as above. We would like to claim as above that a large set of values of x may be distinguished by the same y . However, we have the problem that there are way too many possible values for y for such a pigeonhole argument to work.

Recall that we have constructed a distinguisher A' such that for every $x \in X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m}$, A' $\frac{\epsilon}{2m}$ -distinguishes

$$(C(x)_{y_1}, \dots, C(x)_{y_{i-1}}, b_i, y)$$

from

$$(C(x)_{y_1}, \dots, C(x)_{y_i}, y).$$

Thus, instead of finding a set of x that is distinguished by $y|_{[t]-S_i}$, we can find a set of x that is distinguished by the same functions $(C(x)_{y_1}, \dots, C(x)_{y_{i-1}})$, for any values of $C(x)$. Note that this is not as bad as it appears, since we may define each $C(x)_{y_k}$ by its values for each assignment to

$$y|_{S_i \cap S_k}.$$

Since by definition

$$|S_i \cap S_k| \leq a,$$

we need only define 2^a values of $C(x)_{y_k}$. Thus $(C(x)_{y_1}, \dots, C(x)_{y_{i-1}})$ is determined by $m2^a$ values.

Now we can apply the pigeonhole principle to conclude that there is a distinguisher A'' and a set $X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m, \{C(x)_{y_{<i}}\}}$ of size at least $\frac{\epsilon}{2} 2^{k-m-m2^a}$ on which A'' $\frac{\epsilon}{2m}$ -distinguishes $(C(x)_{y_i}, y_i)$ from (b_i, y_i) for random b_i, y_i .

Recall that by definition, an ϵ -distinguisher only exists if the statistical difference between the quantities being distinguished is at least ϵ . Further, note that y_i represents a uniformly random integer in $[2^k]$, so the statistical difference

$$\|(C(x)_{y_i}, y_i) - (b_i, y_i)\|$$

equals

$$\mathbb{E}_{j \in_R [2^k]} \left[2 \left| \mathbb{E}_x [C(x)_j - \frac{1}{2}] \right| \right].$$

If we define the word d to take for each j the value of the majority of the values in $C(x)_j$ then we can rewrite the statistical difference equivalently as

$$\|(C(x)_{y_i}, y_i) - (b_i, y_i)\| = \frac{1}{2} - \frac{1}{2^k} \mathbb{E}_x [|d - C(x)|],$$

just $\frac{1}{2}$ minus the expected Hamming distance of words in $X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m, \{C(x)_{y_{<i}}\}}$ from d .

Since the ϵ -distinguisher is not given our choice of x , and further, can distinguish $(C(x)_{y_i}, y_i)$ from random for any $x \in X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m, \{C(x)_{y_{<i}}\}}$, we have that all words in $X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m, \{C(x)_{y_{<i}}\}}$ must lie within distance $\frac{1}{2} - \frac{\epsilon}{2m}$ of d .

However, provided $\frac{\epsilon}{2m} > \delta$ where C was defined to correct a $\frac{1}{2} - \delta$ fraction of errors, there can be at most L codewords that lie within distance $\frac{1}{2} - \frac{\epsilon}{2m}$ of d , and we have the desired contradiction provided that

$$|X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m, \{C(x)_{y_{<i}}\}}| \leq L.$$

Above, we have the bound

$$|X_{\frac{\epsilon}{2}, i, b_{i+1}, \dots, b_m, \{C(x)_{y_{<i}}\}}| \geq \frac{\epsilon}{2} 2^{k-m-m2^a},$$

from which we conclude that if

$$k > m2^a + \log O\left(\left(\frac{m}{\epsilon}\right)\right)$$

then we have a (k, ϵ) extractor, as desired. ■

We now compute some parameters of the above construction.

Suppose we wish to construct a (k, ϵ) extractor $\{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$ where $k = \sqrt{n}$, and $m = k^{1-\gamma}$ for some γ . From above, we can construct extractors for

$$2^a \approx \frac{k}{m} = k^\gamma,$$

implying $a \approx \gamma \log k$. Also, for our error-correcting code, $L \approx N = 2^l$, so $l \approx \log N$. Recall that we could construct designs for

$$t \geq \frac{l^2}{a} \exp\left(\frac{\log m}{a}\right).$$

Thus, we have

$$\begin{aligned} t &\approx \frac{l^2}{a} \exp\left(\frac{\log m}{a}\right) \\ &\approx \frac{\log^2 N}{\gamma \log k} \exp\left(\frac{\log k}{\gamma \log k}\right) \\ &\approx \frac{O(\log^2 k)}{\gamma \log k} O(1) = O(\log k). \end{aligned}$$

Thus we use only logarithmic additional randomness. We note that we have seen an existence argument for an extractor that uses

$$t = \log_2 n + O(1)$$

additional randomness, while Trevisan's extractor is a construction requiring only

$$t = O(\log n)$$

additional randomness.

2 The Ta-Shma, Zuckerman, Safra extractor

Let n be the length of the quasi-random source string, and let $q \approx \sqrt{n}$ be prime. We will use a two level construction for the extractor, with a small inner code $C_{\text{small}} : \mathbb{F}_q \rightarrow \{0, 1\}^l$ constructed so as to be $(\frac{1}{2} - \delta, L)$ list decodable used to encode elements in the field \mathbb{F}_q .

The “outer” component of the construction is to consider the length- n input string as specifying the coefficients of a bivariate polynomial P over \mathbb{F}_q of degree (\sqrt{n}, \sqrt{n}) . For now, it does not matter whether the input string is only 0 – 1. The extractor uses the additional randomness to generate an ordered pair $(a, b) \in [\sqrt{n}] \times [\sqrt{n}]$ and an index $j \in [l]$. The output of the extractor is

$$E(P, (a, b), j) = C_{\text{small}}(P(a, b))_j, \dots, C_{\text{small}}(P(a + m - 1, b))_j,$$

where m is the length of the desired output.

We present a heuristic sketch of some of the properties of this extractor. This is covered in more detail in the next lecture.

As for the Trevisan extractor, we suppose for the sake of contradiction that there exists a function A that ϵ -distinguishes the output of this extractor for P drawn randomly from some X of size at least 2^k . As above, by Markov arguments we conclude that A is an $\frac{\epsilon}{2}$ -distinguisher for each $x \in X_{\frac{\epsilon}{2}}$ for some $X_{\frac{\epsilon}{2}}$ of size at least $\frac{\epsilon}{2}2^k$.

Using hybridization arguments similar to those used in the case for the Trevisan extractor, we conclude that there exists a predictor B for $P(a, b)$ given $(P(a - i, b), \dots, P(a - 1, b), a, b)$ where $i < m$. Specifically,

$$B(P(a - i, b), \dots, P(a - 1, b), a, b) = P(a, b)$$

with probability ϵ' over a, b , for every $P \in X'$ where $|X'|$ is large. This will provide the desired contradiction.

Essentially, we show that given the ability to predict values of P , in addition to the fact the P is of low degree in any linear combination of x and y , we can “list-decode” P given only a small amount of extra “free” values of P . We then apply the pidgeonhole principle to show that there is a *specific* set of free values for which we have a large fraction of the polynomials in X' still feasible, which contradicts the fact that we can list-decode from these free values to a *small* set of possible polynomials.

We describe briefly how the fact that P is a bivariate polynomial helps us to form a small list of possible polynomials given a few free values. As noted above, the fact that P is of degree at most \sqrt{n} in x and y implies that the degree of P along any line in \mathbb{F}_q^2 is also at most \sqrt{n} . We now have two orthogonal prediction methods: we can predict $P(a, b)$ given the previous i values on its row, and we can list-decode P on a whole line given some fraction of its values on that line. Our technique will be to evaluate P on i parallel lines, then use our predictor B to predict P on the following line with high probability, and repeat until we have P on every line. One can show by a Chebyshev argument that for some choice of initial line, the set of points for which the predictor B works is distributed fairly evenly along the lines. Thus we can list-decode P on the whole region, yielding the desired contradiction.