# PCP – Lecture 8

Our next goal in the course is to prove

$$NP \subseteq PCP_P[poly, \boxed{1}]$$

$\uparrow$ proximity

← Note the constant query + binary alphabet

We will be able to use this verifier in a composition scheme to prove the PCP thm.

First however, we show a verifier not for NP but for a simpler language — consisting of all linear functions in $n$ variables over $\{0,1\}$

## Linearity Testing

Given $f : \{0,1\}^n \longrightarrow \{0,1\}$

We wish to test whether $f$ is linear. Making <u>few</u> queries to $f$.

**Def:** $f : \{0,1\}^n \longrightarrow \{0,1\}$ is linear if
$$\forall x, y \in \{0,1\}^n \quad f(x) + f(y) = f(x+y) \quad \text{(addition modulo 2)}$$
(of course, this def extends to any group $G$ replacing $\{0,1\}$.)

**Claim:** $f$ is linear iff $\exists a \in \{0,1\}^n$ s.t. $f(x) = \langle a, x \rangle = \sum_i a_i x_i \mod 2$.

**Proof:** ($\Leftarrow$) clear. For ($\Rightarrow$): let $a_i = f(e_i)$ for $e_i = (0 \cdots 0, 1, 0 \cdots 0)$.

and the claim follows by linearity.

Clearly, # queries is at least 3. We now show a 3-query test:

BLR Test : Choose random $x, y \in \{0,1\}^h$.
Test $f(x) + f(y) = f(x+y)$

Clearly if $f$ linear $\implies$ Pr(success) = $\underline{1}$.

What happens if $f$ is not linear? Need to talk about distance.
Def: Let $f, g : \{0,1\}^h \to \{0,1\}$, denote $dist(f,g) = Pr_x \left[ f(x) \neq g(x) \right]$
$f, g$ are $\delta$-far if $dist(f,g) \geq \delta$, $\delta$-close if $dist(f,g) \leq \delta$
$f$ is $\delta$-far from linear if it is $\delta$-far from all linear $g$.

Theorem : If $f$ is $\delta$-far from linear,

$$\text{then } Pr\left[ T \text{ rejects } f \right] \geq \min\left( \frac{2}{9}, \frac{\delta}{2} \right) \geq \frac{2}{9}\delta$$

Comments: * "special case" of low degree test. Preceded it historically.
* can be extended to groups homomorphism testing.
∃ groups for which $\frac{2}{9}$ is tight!
* For our case, $\{0,1\}$, can prove, via Fourier analysis that
$Pr\left[ T \text{ rejects} \right] \geq \delta$.

· let $f : \mathbb{Z}_9^h \to \mathbb{Z}_9$. $f(u) = 3k$ if $u_1 = 3k, 3k-1, 3k+1$
prove (a) $Pr\left[ T \text{ rej } f \right] = \frac{2}{9}$
(b) $Dist(f, Lin) = 2/3$.

**Proof:** Note $f(x) + f(y) \neq f(x+y)$ iff $x_1 = y_1 = 1 \pmod 3$ or $x_1, y_1 = -1 \pmod 3$

this happens w. prob. $2/9$.

Dist$(f, Lin) = 2/3$ : $\cdots$

## Proof of Theorem: We use a correction to majority argument.

<u>Idea</u>: Define a "corrected" version of $f$, $g$ :

For each $x$ consider $f(y) + f(x+y)$ for all $y$.

Define $g(x) = 1$ if $Pr_y[f(y) + f(x+y) = 1] \geq \frac{1}{2}$, and $g(x) = 0$ OW.

Also let $P_x = \text{Prob}_y[f(y) + f(x+y) = g(x)]$. Clearly $\frac{1}{2} \leq P_x \leq 1$.

<u>Claim 1</u>: $\text{Prob}[T \text{ rejects } f] \geq \frac{1}{2} \cdot \text{dist}(g, f)$

<u>Proof:</u>

$$\text{Prob}[T_{rej}] = \underbrace{Pr[g \neq f]}_{=\delta} \underbrace{Pr[T_{rej} \mid g \neq f]}_{\geq \frac{1}{2}} + Pr[g = f] Pr[T_{rej} \mid g = f]$$

$$\geq \delta/2 \qquad\qquad\qquad \square$$

<u>Claim 2</u>: If $Pr(T_{rej}) < \frac{2}{9}$ then $\forall x \quad P_x \geq \frac{2}{3}$.

<u>Proof:</u> $\text{Prob}_{y,z}[f(y) + f(x+y) = f(z) + f(x+z)] = (P_x)^2 + (1-P_x)^2$

rearranging $= \text{Prob}[\underbrace{f(y) + f(z)}_{= f(y+z)} = \underbrace{f(x+y) + f(x+z)}_{= f(y+z)}] > \frac{5}{9}$

$\underbrace{}_{\substack{= f(y+z) \\ \text{w. prob} \geq \frac{7}{9}}}$ $\qquad \underbrace{}_{\substack{= f(y+z) \text{ w. prob} \\ \geq \frac{7}{9}}}$

so $(P_x)^2 + (1-P_x)^2 > \frac{5}{9} \implies P_x > \frac{2}{3}$ $\qquad \square$

Claim 3: $g$ is linear.

Proof:

$$g(x) \; + \; g(y) \qquad\qquad g(x+y)$$

$$f(z) \quad f(x+z) \quad f(z) \quad f(z+y) \qquad f(z+x) \quad f(z-y)$$

$$\text{for} > \tfrac{2}{3} \; z\text{'s} \qquad \text{for} > \tfrac{2}{3} \; z\text{'s} \qquad \text{for} > \tfrac{2}{3} \; z\text{'s}$$

$\exists z^*$ for which all three hold. $\Rightarrow$

$$\underbrace{f(z^*) + f(x+z^*)}_{g(x)} + \underbrace{f(z^*) + f(y+z^*)}_{g(y)} = \underbrace{f(x+z^*) + f(z^*+y)}_{g(x+y)}$$

$$g(x) \qquad g(y) \qquad = \qquad g(x+y) \qquad \square$$

Conclusion: $\text{Prob}\left[ T \text{ rejects} \right]$ is either $\geq \tfrac{2}{9}$ or

$$\text{Prob}\left[ T_{\text{rej}} \right] \geq \frac{\delta}{2} = \frac{\text{dist}(f,g)}{2} \geq \frac{\text{dist}(f, \text{Lin})}{2}.$$

$g$ linear

Def: The Hadamard Code maps to each $a \in \{0,1\}^h$ the linear function $L_a : \{0,1\}^h \to \{0,1\}$ defined by $h_a(x) = \sum_i a_i x_i \pmod 2$.

It is an error correcting code $H: \{0,1\}^h \to \{0,1\}^{2^h}$ such that

ⓐ $\forall a \neq b \quad \text{dist}(h_a, h_b) = \tfrac{1}{2} \cdot 2^h$. ( relative distance $= \tfrac{1}{2}$ )

ⓑ Its rate is logarithmic ( $\log N$ bits are mapped to $N$ bits)

ⓒ It is locally testable with 3 queries.

it is a  Locally Testable Code.

**Thm #2:** $\text{Prob}(T \text{ rejects } f) \geq \text{dist}(f, \text{Lin})$

**Proof #2:** (based on Fourier Analysis)

[ we switch notation $0 \to 1$ $1 \to -1$ $a \to (-1)^a$ ]

So a linear function is now $f(x) \cdot f(y) = f(xy)$ pointwise

Fix a function $f : \{\pm 1\}^n \to \{\pm 1\} \subseteq \mathbb{R}$. The space of all functions $f \in \mathbb{R}^{2^n}$ is a vector space.

The standard basis is $\{e_w\}_{w \in \{\pm 1\}^n}$ : $e_w(x) = \begin{cases} 1 & x = w \\ 0 & \text{otherwise} \end{cases}$

Another basis is the following

$\forall \, S \subseteq [n]$ let $\chi_S(x_1 \dots x_n) = \prod_{i \in S} x_i$. Clearly $\chi_S : \{\pm 1\}^n \to \{\pm 1\}$.

Define an inner product $\langle f, g \rangle = 2^{-n} \sum_x f(x) g(x)$.

ⓐ $\langle \chi_S, \chi_S \rangle = 1$

ⓑ $S \neq T$ $\langle \chi_S, \chi_T \rangle = 2^{-n} \sum_x \prod_{i \in S} x_i \prod_{i \in T} x_i$

$$= 2^{-n} \sum_x \prod_{S \Delta T} x_i = 0$$

by pairing $x$ according to their val on some coor in $S \Delta T$.

So $\{\chi_S\}$ is a basis, and $\forall f$ $f = \sum_S \underbrace{\langle f, \chi_S \rangle}_{\text{notation: } \hat{f}_S} \cdot \chi_S$

ⓒ $\chi_S$ is linear : $\chi_S(x) \cdot \chi_S(y) = \prod_{i \in S} x_i y_i = \chi_S(xy)$.

(d) for any $f: \{\pm 1\}^h \to \{\pm 1\}$

$$\hat{f}_s = \langle f, \chi_s \rangle = Pr(f = \chi_s) - Pr(f \neq \chi_s) = 1 - 2\,dist(f, \chi_s).$$

(e) for any $f: \{\pm 1\}^h \to \{\pm 1\}$     $\sum (\hat{f}_s)^2 = 1$

$$\langle f, f \rangle = 2^{-n} \sum_x (f(x))^2 = 1$$

$$\langle f, f \rangle = \langle \sum_s \hat{f}_s \chi_s, \sum_s \hat{f}_s \chi_s \rangle = \sum \hat{f}_s^2$$

Now we want to relate $S = \underset{x,y}{Prob}\left( f(x)f(y) \neq f(xy) \right) = \underset{x,y}{Prob}\left( f(x)f(y)f(xy) \neq 1 \right)$ to the F. coef. of $f$.

Let $e = \underset{x,y}{\mathbb{E}}\left( f(x)f(y)f(xy) \right)$ then $e = S - (1-S) = 1 - 2S$.

$$e = \underset{xy}{\mathbb{E}}\left[ \sum_s \hat{f}_s \chi_s(x) \sum_T \hat{f}_T \chi_T(y) \sum_U \hat{f}_U \chi_U(xy) \right]$$

$$= \underset{xy}{\mathbb{E}} \sum_{STU} \hat{f}_s \hat{f}_T \hat{f}_U \prod_{i \in S} x_i \prod_{i \in T} y_i \prod_{i \in U} x_i y_i$$

$$= \underset{xy}{\mathbb{E}} \sum_{STU} \hat{f}_s \hat{f}_T \hat{f}_U \prod_{i \in S \triangle U} x_i \prod_{i \in T \triangle U} y_i = \sum_s (\hat{f}_s)^3$$

$$\leq \max_s \hat{f}_s \cdot \underbrace{\sum \hat{f}_s^2}_{=1} = \max_s \hat{f}_s = \hat{f}_{s_0} \quad \left(\text{denoting } s_0 \text{ a maximal coef}\right)$$

$$\Rightarrow S = \frac{1-e}{2} \geq \frac{1 - \hat{f}_{s_0}}{2} = \frac{1 - (1 - 2\,dist(f, \chi_{s_0}))}{2} = dist(f, \chi_{s_0})$$

$$= dist(f, Lin)$$

This completes our analysis of the BLR linearity testing.

# Self - Correction

As in the low degree case, the Had code allows self correction.

**Lemma:** $f: \{0,1\}^n \longrightarrow \{0,1\}$, $\text{dist}(f, \text{Lin}) \leq \delta < \frac{1}{4}$.

Then ⓐ There is a unique linear $g: \{0,1\}^n \longrightarrow \{0,1\}$ that is closest to $f$.

ⓑ There is a randomized two-query procedure $S$ that guarantees for every $x$: $\Pr_r[S(x) = g(x)] \geq 1 - 2\delta > \frac{1}{2}$.

**Proof:** ⓐ If there were $g_1, g_2$ linear, both $\delta$-close to $f$ then (by $\triangle$ ineq.:)

$$\frac{1}{2} \leq \text{dist}(g_1, g_2) \leq \text{dist}(g_1, f) + \text{dist}(f, g_2) \leq 2\delta < \frac{1}{2}$$

contradiction.

ⓑ $S(x)$: choose random $y \in \{0,1\}^n$, output
$$f(y) + f(x+y).$$

Since $y$ is unif. dist. it hits the set $\text{BAD} = \{x \mid g(x) \neq f(x)\}$ with prob. $\leq \delta$, and similarly $x+y$. Altogether:

$$\text{Prob}\left(f(y) + f(x+y) = g(x)\right) \geq \text{Prob}\left(y \in B \text{ OR } x+y \in B\right) \geq 1 - 2\delta > \frac{1}{2}.$$

... good for program checking.

**Def:** An ecc $C: \{0,1\}^k \longrightarrow \{0,1\}^n$ is <u>locally decodable</u> with $q$-queries if there is a randomized procedure $D$ such that, given $w \in \{0,1\}$, $\text{dist}(w, C) \leq \delta$, on input $i$ $D$ outputs $x_i$ where $C(x)$ is the codeword closest to $w$.
s.t. $D$ makes only $\leq q$ queries into $w$.

**Claim:** The Hadamard Code is locally decodable with $2$ queries.

<u>Rmk</u>: On input $i$ $D$ runs $S(e_i)$ where $e_i = (0 \text{---} 0, \overset{i}{\underset{\downarrow}{1}}, 0 \text{---} 0)$.

<span style="color:magenta"><u>research question</u>:</span> <span style="color:magenta">are there locally decodable codes with good rate?</span>
<span style="color:magenta">(polynomial?)</span> <span style="color:blue">with 2-queries ← no!</span>
<span style="color:blue">3-queries ???</span>

<span style="color:magenta">this is also a "popular" lower bound question</span>

In fact, the Hadamard code allows one to read "correctly"
not only the value of $x_i$, but also the value of $l(x)$ for
<u>all</u> linear functions $l$ (by the "self correction" property).

We will now slightly strenthen this property to all $q(x)$ for all
quadratic functions $q(x_1 \ldots x_n) = \sum a_{ij} x_i x_j + \sum l_i x_i + l_0$.

<u>Def</u>: The quadratic functions encoding $Q: \{0,1\}^n \longrightarrow \{0,1\}^{n^2}$
maps $(a_2 \ldots a_n) \in \{0,1\}^{n-1}$ into $H(a' \otimes a')$
where $a' \otimes a' \in \{0,1\}^{n^2}$ is defined by $(a' \otimes a')_{ij} = a'_i \cdot a'_j$
and $a' \in \{0,1\}^n$ is the vector $(1, a_2 \ldots a_n)$.

― The distance of this code is $\geq \frac{1}{4}$.
― The rate is $\sqrt[4]{n} \ldots$
― locally testable? locally decodable?

**Claim 1:** $Q$ is locally decodable. Moreover, for any quadratic function $g(x_1 ... x_n)$ there is a two-query procedure that whp gives $g(x_1 ... x_n)$ if given oracle access to $f \in \{0,1\}^{n^2}$ s.t. $\text{dist}(f, Q) \leq \delta < \frac{1}{4}$.

**Proof:** Since $\text{Im}(Q)$ is a subset of $\text{Im}(Had)$ the local decoding procedure for Had works for $Q$. Moreover, every quadratic function $q$ on $\vec{a}$ can be expressed as a linear function on $b = a' \otimes a'$. So by the self-correction (strong) of $Q$ we get the result.

**Claim 2:** $Q$ is locally Testable.