

Learning using Local Membership Queries

Pranjal Awasthi(CMU)

Vitaly Feldman(IBM Almaden)

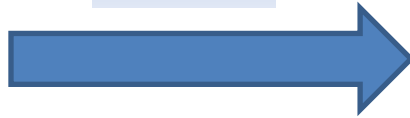
Varun Kanade(UC Berkeley)

PAC + MQ Learning [Angluin '87]

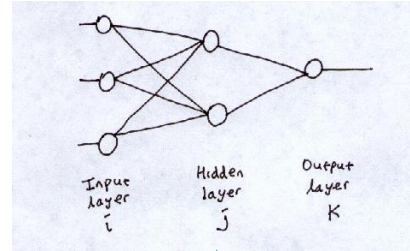
Distribution Oracle



$x \gg D$
 $(x, f(x))$

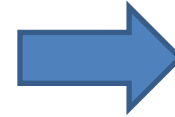


Learning algorithm



x $f(x)$

f



$h: \Pr[f \neq h] \cdot 2$

PAC + MQ Learning

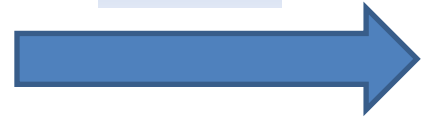
- Can learn decision trees [Kushilevitz-Mansour '91, Bshouty '93]
 - Agnostic learning [Kalai et al. '08]
- DNFs under product distributions [Jackson '94]
- Learning intersection of 2 halfspaces [Baum '91, Gopalan et al. '12]

PAC + MQ Learning [Angluin '87]

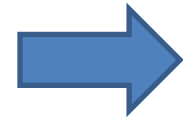
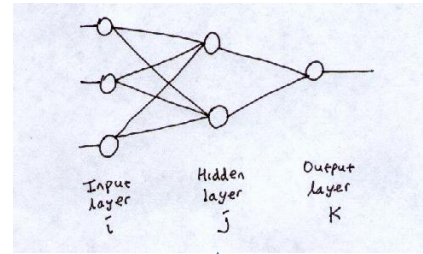
Distribution Oracle



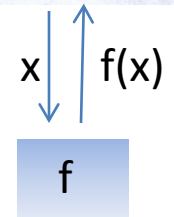
$x \gg D$
 $(x, f(x))$



Learning algorithm



$h: \Pr[f \neq h] \cdot 2$



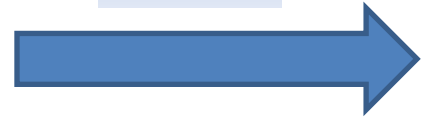
Too Strong

PAC + MQ Learning [Angluin '87]

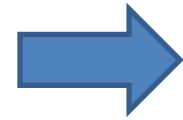
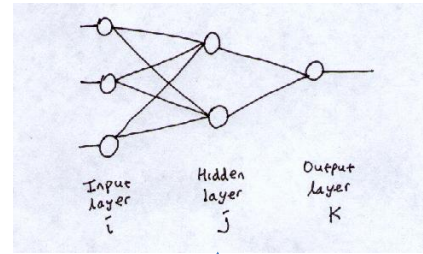
Distribution Oracle



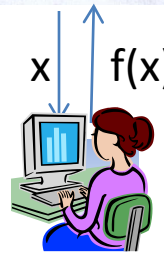
$x \gg D$
 $(x, f(x))$



Learning algorithm



$h: \Pr[f \neq h] \cdot 2$



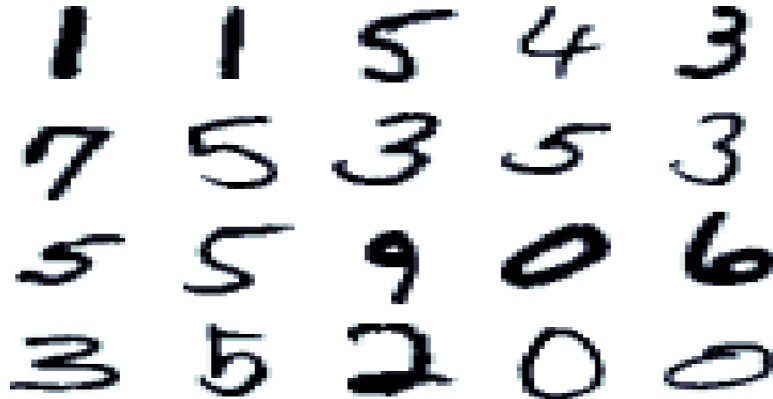
Too Strong

Not every x will make sense to a labeler

PAC + MQ learning

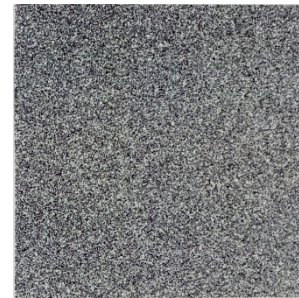
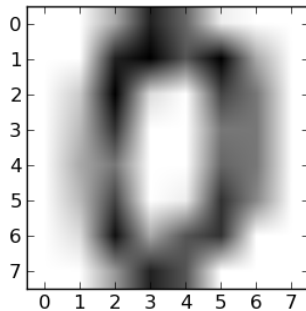
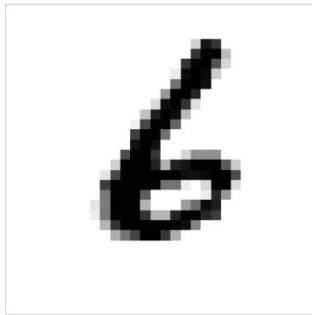
[Lang, Baum'93]:

Query algorithms do not work well when a human oracle is used.



Our model

1 1 5 4 3
7 5 3 5 3
9 5 9 0 6
3 5 2 0 0

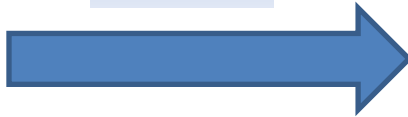


PAC + local-MQ Learning

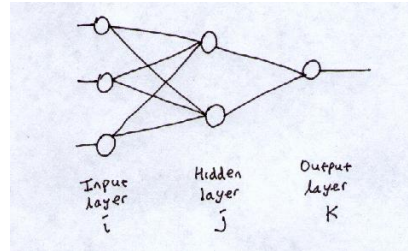
Distribution Oracle



$x \gg D$
 $(x, f(x))$



Learning algorithm



$h: \Pr[f \neq h] \cdot 2$

Local query



$x \updownarrow f(x)$

f

A query x is **r-local** if there exists a data point x' in the sample such that $d(x, x') \leq r$.

For this talk:
 $x \in \{0, 1\}^n$
 $d = \text{Hamming}$

Warmup: 1-local MQs

Claim:

- Can learn noisy parities using 1-local MQs.
 - Compute $\Pr[f(x) \text{ changes} \mid x_i \text{ is flipped}]$
- Can learn k -juntas under uniform distribution.

Our Results(1)

Theorem:

Efficient $O(\log(n))$ -local algorithm for learning sparse polynomials over $\{0,1\}^n$ under \log -Lipschitz distributions.

Includes log-depth decision trees.

Our Results(2)

Theorem:

Efficient $O(\log(n))$ -local algorithm for learning poly size decision trees under constant biased **product** distributions.

Our Results(3)

Theorem:

$O(\log(n))$ -local algorithm for learning poly size DNFs under **uniform** distribution.

Run time is $n^{O(\log\log(n))}$

Learning sparse polynomials

Theorem:

Efficient $O(\log(n))$ -local algorithm for learning sparse polynomials over $\{0,1\}^n$ under \log -Lipschitz distributions.

Learning sparse polynomials

- $X = \{0,1\}^n$
- D is \mathbb{R} -log-Lipschitz
 - $\exists x, x'$ s.t. $d(x, x') = 1$, $D(x)/D(x') \leq \mathbb{R}$
 - $\mathbb{R} = 1$, uniform
 - Includes constant biased product distributions, smooth distributions of [Kalai et al. '09].

Learning sparse polynomials

- $X = \{0,1\}^n$
- Target function: $f(x) = \sum_{i \in S} c_i x_i$
- Bounded coefficients
- $|\text{supp}(f)| = \text{poly}(n)$
- Goal: Output h s.t. $\|f - h\|_2 \leq \epsilon$

Learning sparse polynomials

- Low degree approximation
- Identifying candidate low degree monomials
- L_2 regression over candidate monomials.

Learning sparse polynomials

- **Low degree approximation**
- Identifying candidate low degree monomials
- L_2 regression over candidate monomials.

Low degree approximation

- $f(x) = \sum_{S: |S|=d} c_S \prod_{i \in S} x_i$
- $f_d(x) = \sum_{S: |S| \leq d} c_S \prod_{i \in S} x_i$; $d = O(\log(n^2))$
- **Claim:** $\|f - f_d\|_2 \leq \epsilon$

Proof:

- For any S , $\left(\frac{1}{1+\epsilon}\right)^{|S|} \cdot \Pr_D \left[\prod_{i \in S} x_i = 1 \right] \cdot \left(\frac{\epsilon}{1+\epsilon}\right)^{|S|}$

Learning sparse polynomials

- Low degree approximation
- Identifying candidate low degree monomials
- L_2 regression over candidate monomials.

Kushilevitz-Mansour Algorithm

Theorem[KM '91]:

For any $f: \{0,1\}^n \rightarrow \{0,1\}$, can output all μ heavy monomials in time $\text{poly}(n, 1/\mu, L_1(f))$

Kushilevitz-Mansour Algorithm

L₂ test

Is $E[(f_1)^2]$ large?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

f_1

$$x_1 + x_1 x_2 + x_3 x_1 + x_1 x_2 x_3$$

f_{-1}

$$x_2 + x_3 + x_2 x_3$$

Kushilevitz-Mansour Algorithm

L₂ test

Is $E[(f_1)^2]$ large?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

f_1

$$x_1 + x_1 x_2 + x_3 x_1 + x_1 x_2 x_3$$

f_{-1}

$$x_2 + x_3 + x_2 x_3$$

$f_{1,2}$ x_2

$$x_1 x_2 + x_1 x_2 x_3$$

$f_{1,-2}$ x_2

$$x_1 + x_3 x_1$$

$f_{-1,2}$ x_2

$$x_2 + x_2 x_3$$

$f_{-1,-2}$ x_2

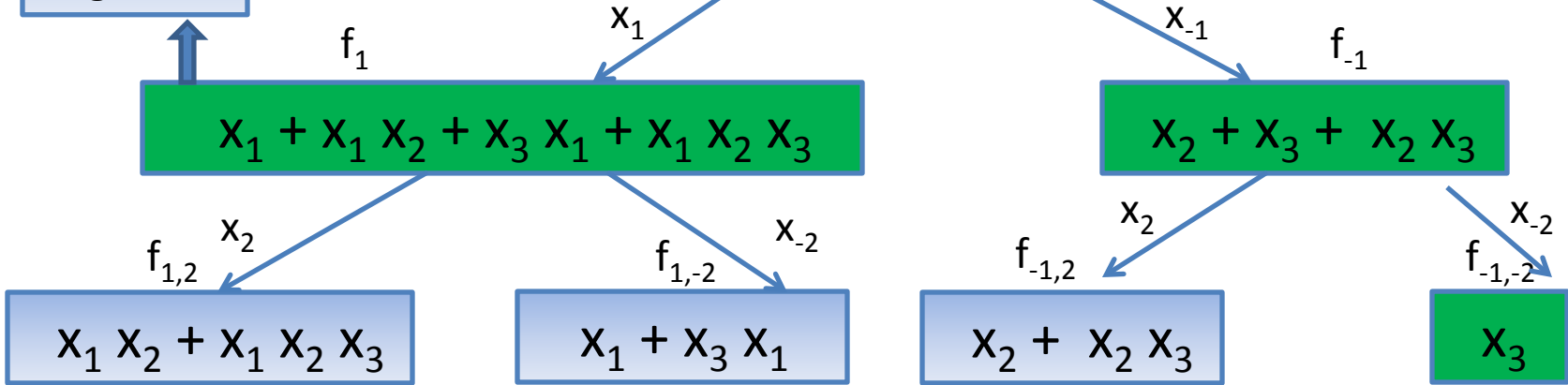
$$x_3$$

Kushilevitz-Mansour Algorithm

L₂ test

Is $E[(f_1)^2]$ large?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$



Leaves represent important monomials.

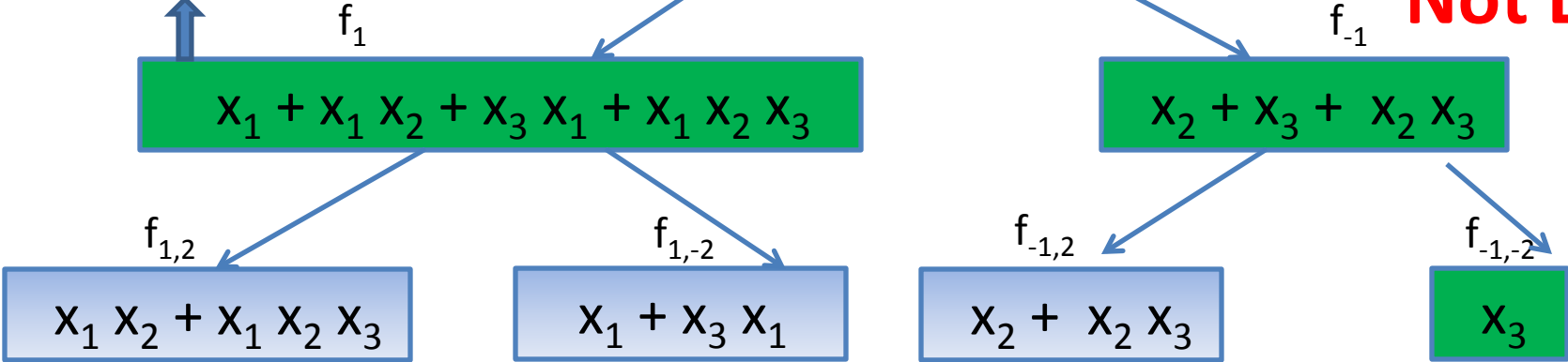
Kushilevitz-Mansour Algorithm

L₂ test

Is $E[(f_1)^2]$ large?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

Not Local



At each step “active” nodes partition the polynomial.

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$
noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

$$x_1 + x_1 x_2 + x_3 x_1 + x_1 x_2 x_3$$

$$x_2 + x_1 x_2 + x_2 x_3 + x_1 x_2 x_3$$

$$x_3 + x_1 x_3 + x_2 x_3 + \dots$$

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

f_1

$$x_1 + x_1 x_2 + x_3 x_1 + x_1 x_2 x_3$$

f_2

$$x_2 + x_1 x_2 + x_2 x_3 + x_1 x_2 x_3$$

f_3

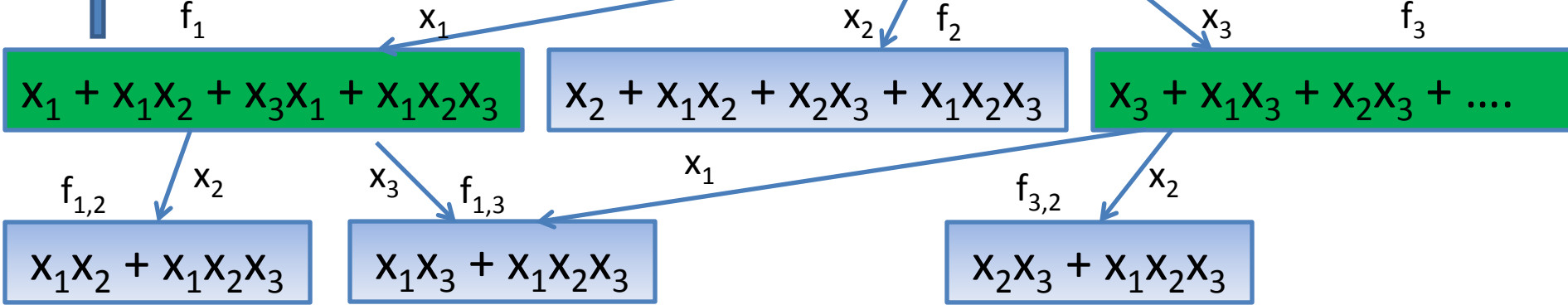
$$x_3 + x_1 x_3 + x_2 x_3 + \dots$$

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

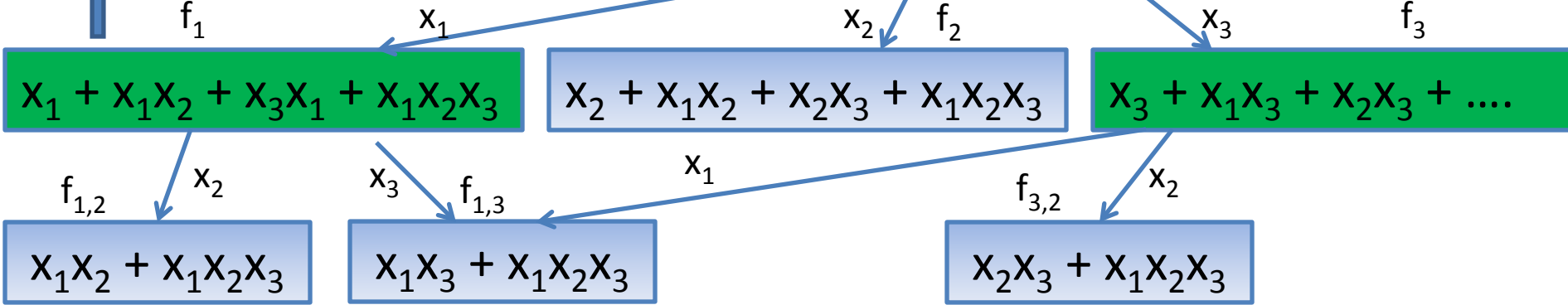


Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$

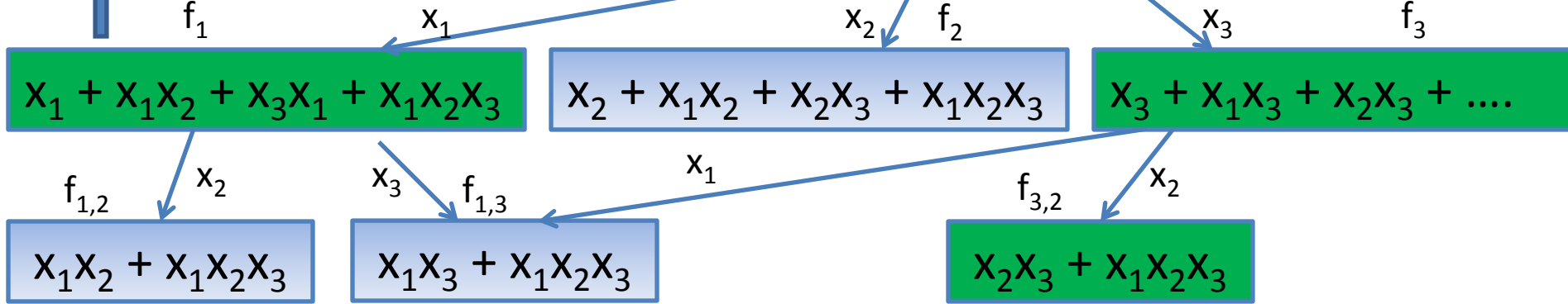


Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$



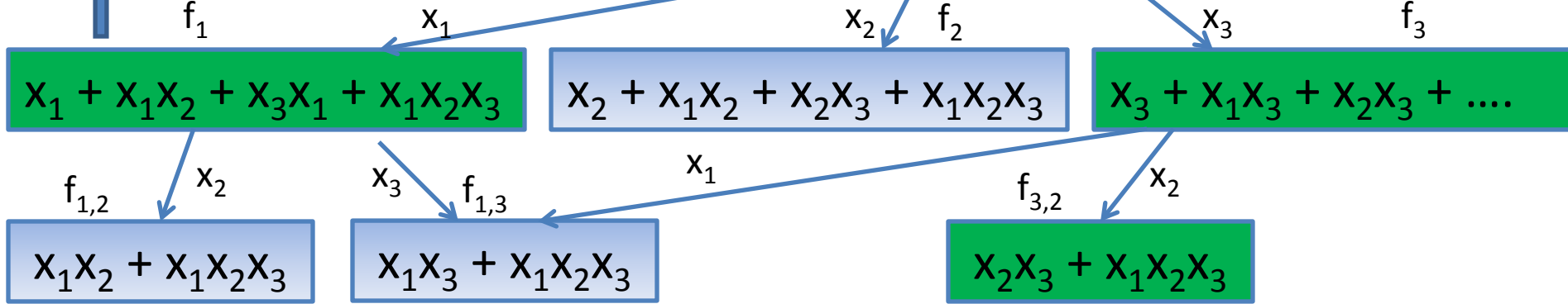
Till depth $d = \log(n/2)$

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$



If f_s passes the test, add S

Non-zero-test(x_1)

$$f(x_1; x_2; \text{ccc}; x_n) = x_1 \overbrace{f_1(x_2; x_3; \text{ccc}; x_n)} + f_{i-1}(x_2; x_3; \text{ccc}; x_n)$$

The test passes if $\Pr_D [f_1 \neq 0]$ is noticeable.

Claim:

Non-zero-test can be implemented using local MQs.

Non-zero-test(x_1)

$$f(x_1; x_2; \dots; x_n) = x_1 f_1(x_2; x_3; \dots; x_n) + f_{i-1}(x_2; x_3; \dots; x_n)$$

$$f(1; x_2; \dots; x_n) = f_1(x_2; x_3; \dots; x_n) + f_{i-1}(x_2; x_3; \dots; x_n)$$



$$f(0; x_2; \dots; x_n) = 0 f_1(x_2; x_3; \dots; x_n) + f_{i-1}(x_2; x_3; \dots; x_n)$$

$$f_1(x_2; \dots; x_n) = \frac{f(1; x_2; \dots; x_n) - f(0; x_2; \dots; x_n)}{1 - 0} = f(1; x_2; \dots; x_n) - f(0; x_2; \dots; x_n)$$

Non-zero-test(S)

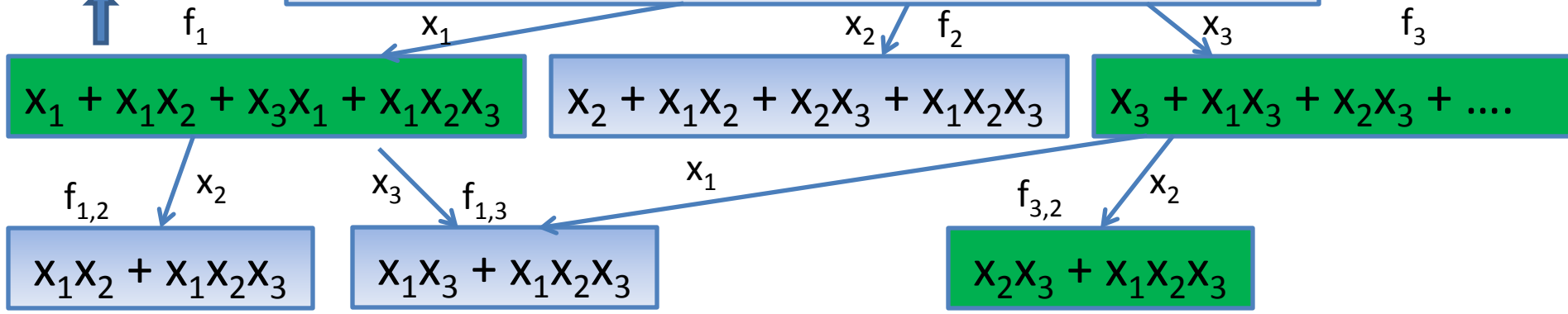
- $f(x) = \sum_{x_i \in S} \tilde{x}_i f_S(x) + f_{i \notin S}(x)$
- $f_S(x) = \sum_{x_i \in S} \tilde{x}_i (2x_i - 1) f(x)$
- Needs $|S|$ -local queries.

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$



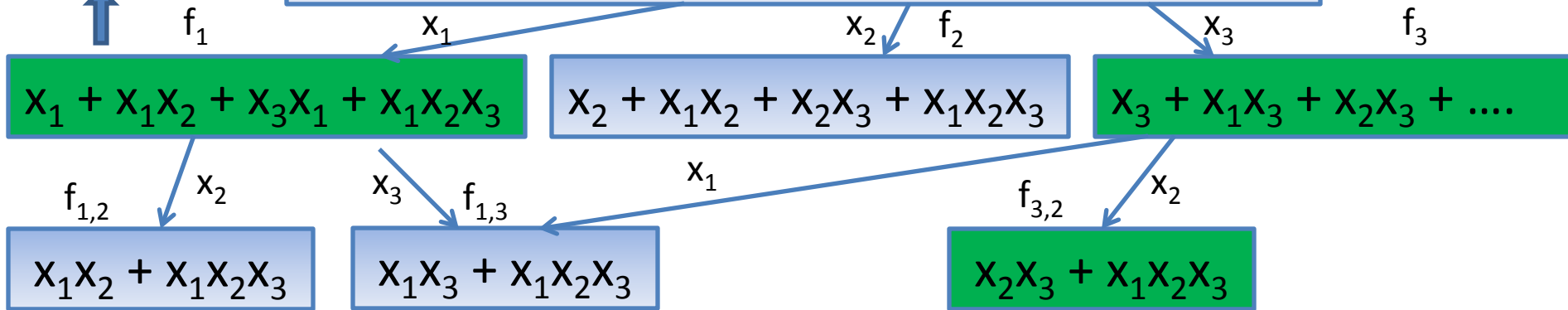
Till depth $d = \log(n/2)$

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$



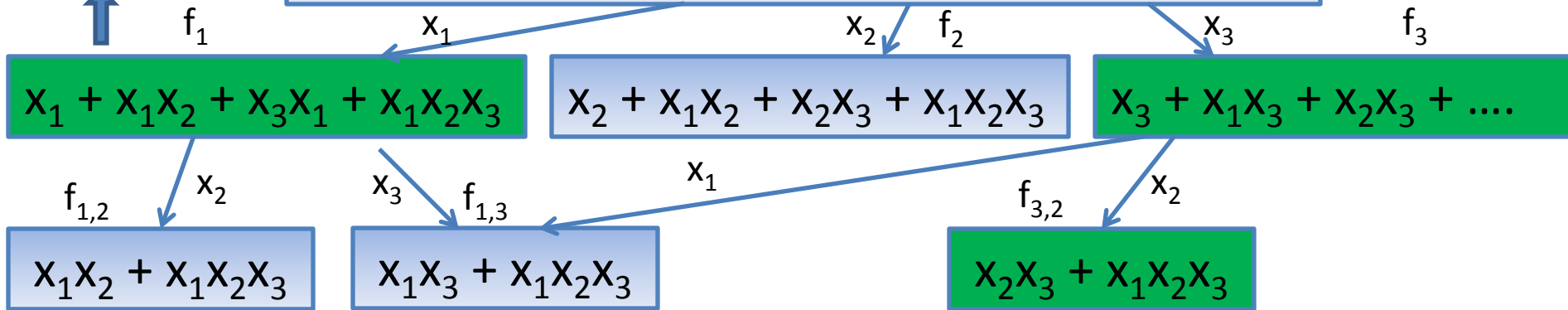
Fit a polynomial $p(x)$ over added terms which minimizes $E((p(x) - f(x))^2)$.

Our Algorithm

Non-zero-test

Is $\Pr[f_1 \neq 0]$ noticeable?

$$f(x) = x_1 + x_2 + x_3 + x_1 x_2 + x_2 x_3 + x_3 x_1 + x_1 x_2 x_3$$



It Works!

Conclusions

log(n)-local algorithms for

- Sparse polynomials under \mathbb{R} log-Lipschitz distributions.
- Decision trees under product distributions.
- DNFs under uniform distribution in time $n^{O(\log \log(n))}$.

Open Problems

- Algorithms for poly size decision trees under log-Lipschitz distributions?
- Polynomial time algorithm for DNFs under uniform distribution?
- Local-MQ agnostic learning algorithms?
 - $O(1)$ -local queries won't help.

THANK YOU

Correctness

- $f(x) = \sum_{S \subseteq [n]} c_S \prod_{i \in S} x_i$
- $f_d(x) = \sum_{S: |S| \leq d} c_S \prod_{i \in S} x_i$; $d = O(\log(n^2))$

Claim:

$\forall S \subseteq [n], \Pr[f_d(x) \neq 0] \geq \exp(-d) = \Omega(1/n^2)$

Correctness

- $f(x) = \sum_{S: |S|=d} c_S \prod_{i \in S} x_i$
- $f_d(x) = \sum_{S: |S|=d} c_S \prod_{i \in S} x_i$; $d = O(\log(n^2))$

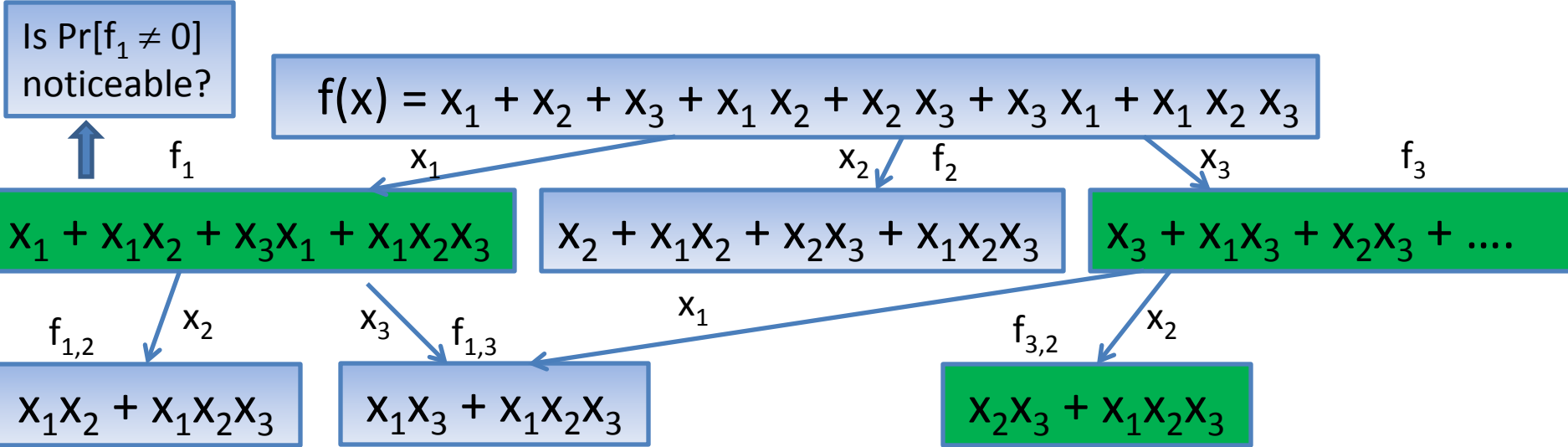
Lemma:

Let f be t -sparse with non-zero constant term and D is

$$\mathbb{R}\text{-log-Lipschitz: } \Pr[f \neq 0] \geq (1 + \mathbb{R})^{-\log(t)}$$

Efficiency

Non-zero-test



Claim:

If S is added then $\exists T \perp S$ s.t. $c_T \neq 0$ and $|T| \cdot d' = O(d)$

Efficiency

Claim:

If S is added then $\exists T \perp S$ s.t.

1. $c_T \neq 0$

2. $|T| \cdot d' = O(d)$



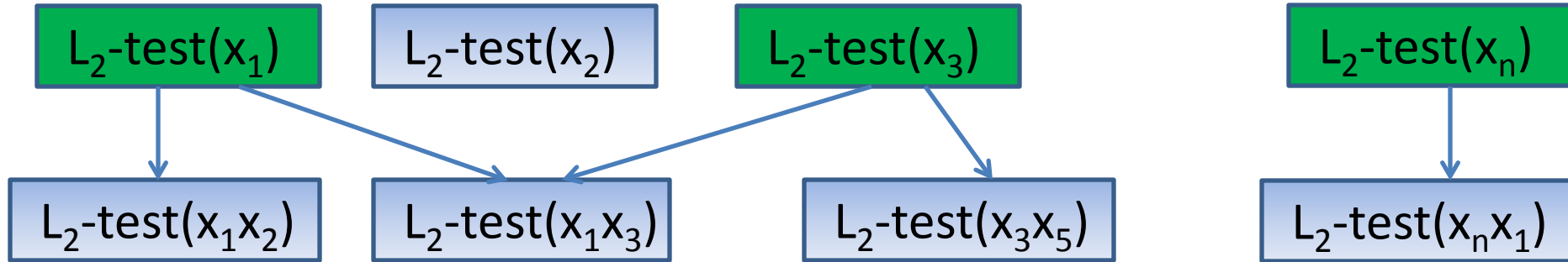
terms added \cdot
 $\text{poly}(n) 2^{O(d)}$

Learning DNFs

Low degree approximation

- $f(x) = \sum_S c_S \tilde{A}_S(x)$
 - $f_d(x) = \sum_{S: |S| \leq d} c_S \tilde{A}_S(x); \quad d = O(\log(\frac{n}{2}))$
 - Finding heavy low-degree coefficients is enough.
-
- Poly-size DNF well approximated by a $O(\log(n))$ -term DNF.
 - For an s -term DNF access to heavy coefficients up to degree $\log(s/2)$ is enough [Feldman '12].

The Algorithm

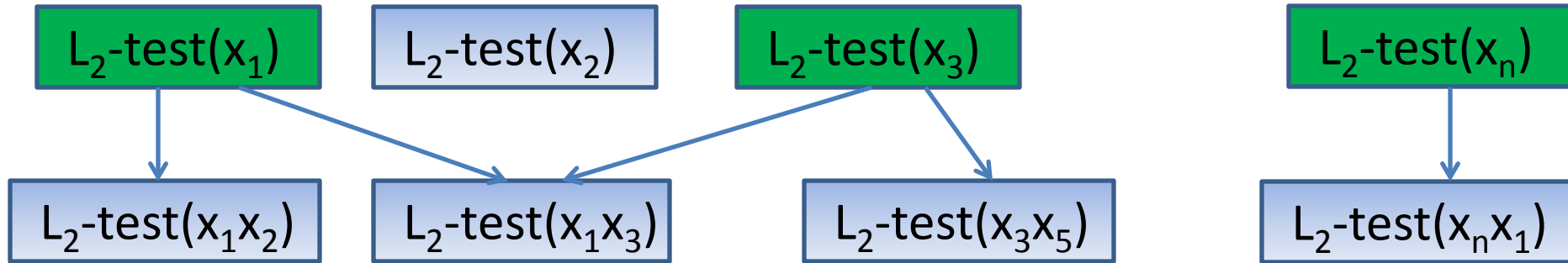


L_2 -test(S)

- $f(x) = \bar{A}_S(x)f_S(x) + f_{i_S}(x)$
- Add S if $E[(f_S(x))^2]$ is large
- Needs $|S|$ -local queries.

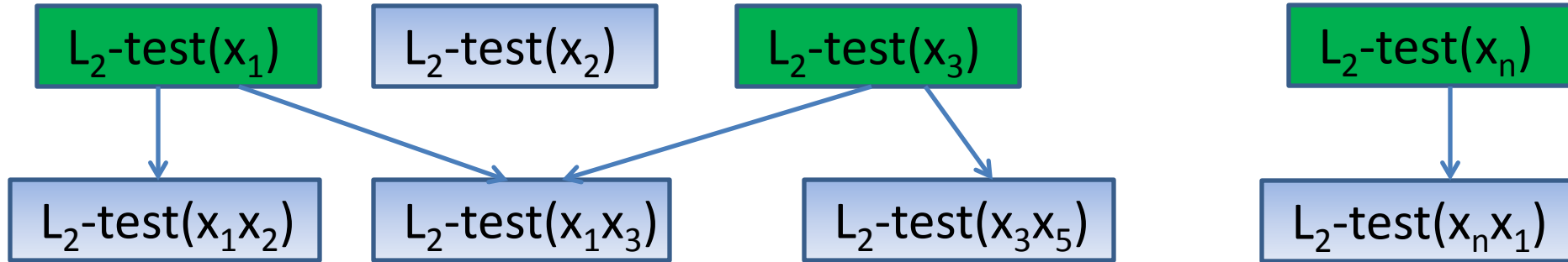
$$8S; E[f_S^2(x)] = \frac{1}{|S|} \sum_{i \in S} \sigma_i^2$$

Correctness



- $f(x) = \sum_S c_S \bar{A}_S(x)$
- $f_d(x) = \sum_{S:j \in S} c_S \bar{A}_S(x)$
- $\sum_S; E[f_S^2(x)] = \sum_{T \parallel S} c_T^2$

Efficiency



- # nodes tested $\cdot n^{O(\log\log(n))}$

Noise in MQ model

- Each query answered incorrectly w.p. ϵ
[Sakakibara '91]
- Persistent noise model [Goldman et al. '90, Angluin, Slonim '91]
- Unreliable boundary queries [Blum et al. '95]

Learning Decision Trees

- $X = \{-1, 1\}^n$, distribution D is uniform.
- $f(x) = \sum_S c_S \hat{A}_S(x); |c_S| \leq \frac{\text{poly}(n)}{2^{|S|}}; \|f\|_1 = \text{poly}(n)$
- Goal: Output h s.t. $\Pr_D [f(x) \neq h(x)] \leq \epsilon$

Learning Decision Trees

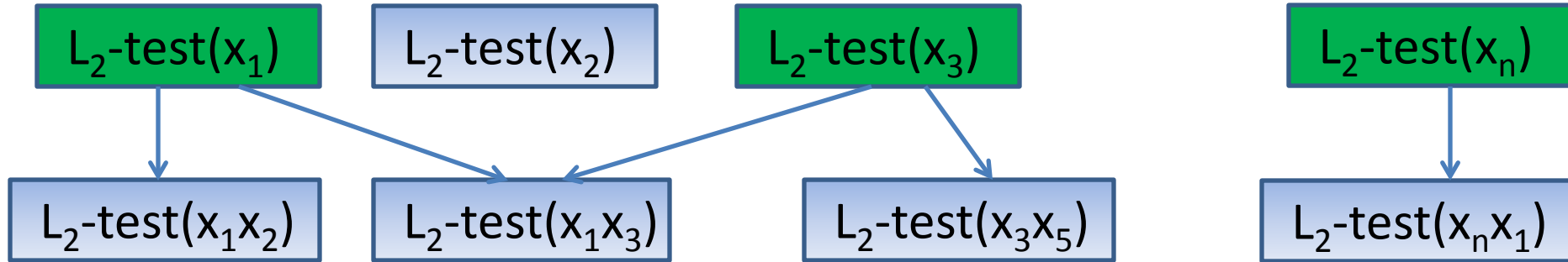
Low degree approximation

- $f(x) = \sum_S c_S \tilde{A}_S(x)$
- $f_d(x) = \sum_{S: |S| \leq d} c_S \tilde{A}_S(x)$
- $\|f - f_d\|_2 \leq \epsilon^2$; $d = O(\log(n/\epsilon^2))$

Proof

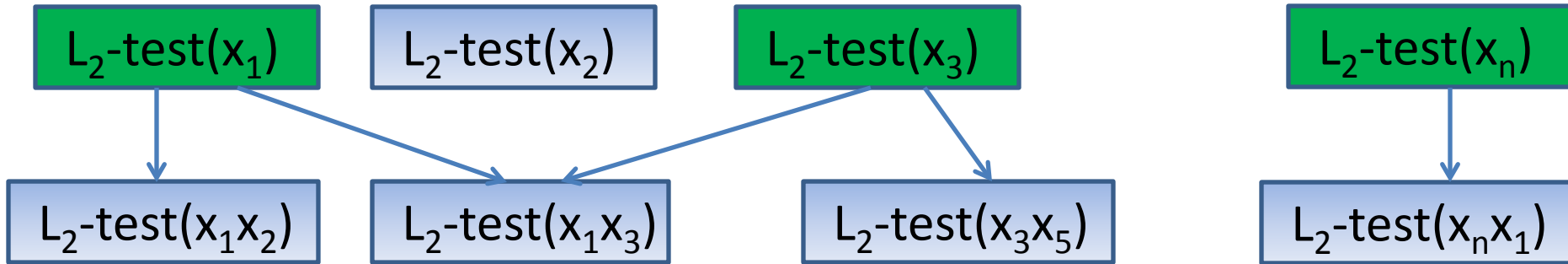
- For any S , $c_S \leq \frac{\text{poly}(n)}{2^{|S|}}$

Correctness



- $f(x) = \sum_S c_S \tilde{A}_S(x)$
- $f_d(x) = \sum_{S:|S|=d} c_S \tilde{A}_S(x)$
- $\sum_S E[f_S^2(x)] = \sum_{T \parallel S} c_T^2$

Efficiency

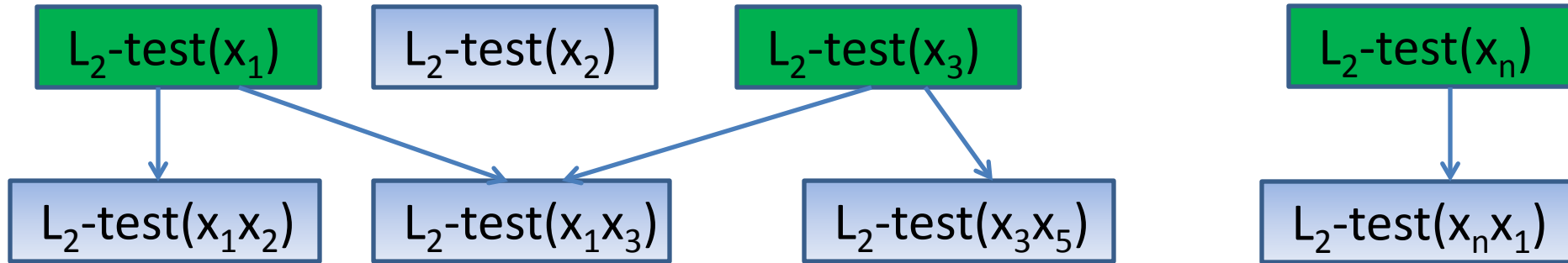


- $f(x) = \sum_S c_S \bar{A}_S(x)$

- $f_d(x) = \sum_{S:|S|=d} c_S \bar{A}_S(x)$

- $\delta_S; E[f_S^2(x)] = \sum_{T \supseteq S} c_T^2 \cdot \max_{T \supseteq S} |c_T|$

Efficiency



- $f(x) = \sum_S c_S \bar{A}_S(x)$
- $f_d(x) = \sum_{S:|S|=d} c_S \bar{A}_S(x)$
- If S is added then exists $T \supset S$, s.t. $|c_T|$ is large.

