

Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals

Clayton Sanford
Columbia Computer Science
February 14th, 2022



Near-Optimal Statistical Query Lower Bounds for **Agnostically Learning** Intersections of Halfspaces with Gaussian Marginals

- **Goal:** Learn $f: \mathbb{R}^n \rightarrow \{\pm 1\}$ from hypothesis class \mathcal{H} given m samples $(x_i, y_i)_{i \in [m]} \sim \mathcal{D}$.
- From m samples, obtain $f \in \mathcal{H}$ satisfying $\Pr[h(x) \neq y] \leq \text{OPT} + 0.1$ for $\text{OPT} = \min_{h \in \mathcal{H}} \Pr[h(x) \neq y]$ with probability ≥ 0.9 .

Near-Optimal **Statistical Query** Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals

- **Problem:** Hard to prove computational lower bounds of agnostic learning.
- **Idea:** Use a restricted model for lower bounds: queries instead of samples.
- A **statistical query** takes (g, τ) with $g : \mathbb{R}^n \times \{\pm 1\} \rightarrow [-1, 1]$, tolerance $\tau > 0$ and returns $q \in \mathbb{E}[g(x, y)] \pm \tau$.
- **Goal:** Learn $f : \mathbb{R}^n \rightarrow \{-1, 1\}$ using SQs.
- For $\text{poly}(n)$ queries of tolerance $\tau \geq 1/\text{poly}(n)$, obtain f satisfying $\Pr[f(x) \neq y] \leq \text{OPT} + 0.1$ for $\text{OPT} = \min_{h \in \mathcal{H}} \Pr[h(x) \neq y]$ with probability ≥ 0.9 .

Near-Optimal **Statistical Query** Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals

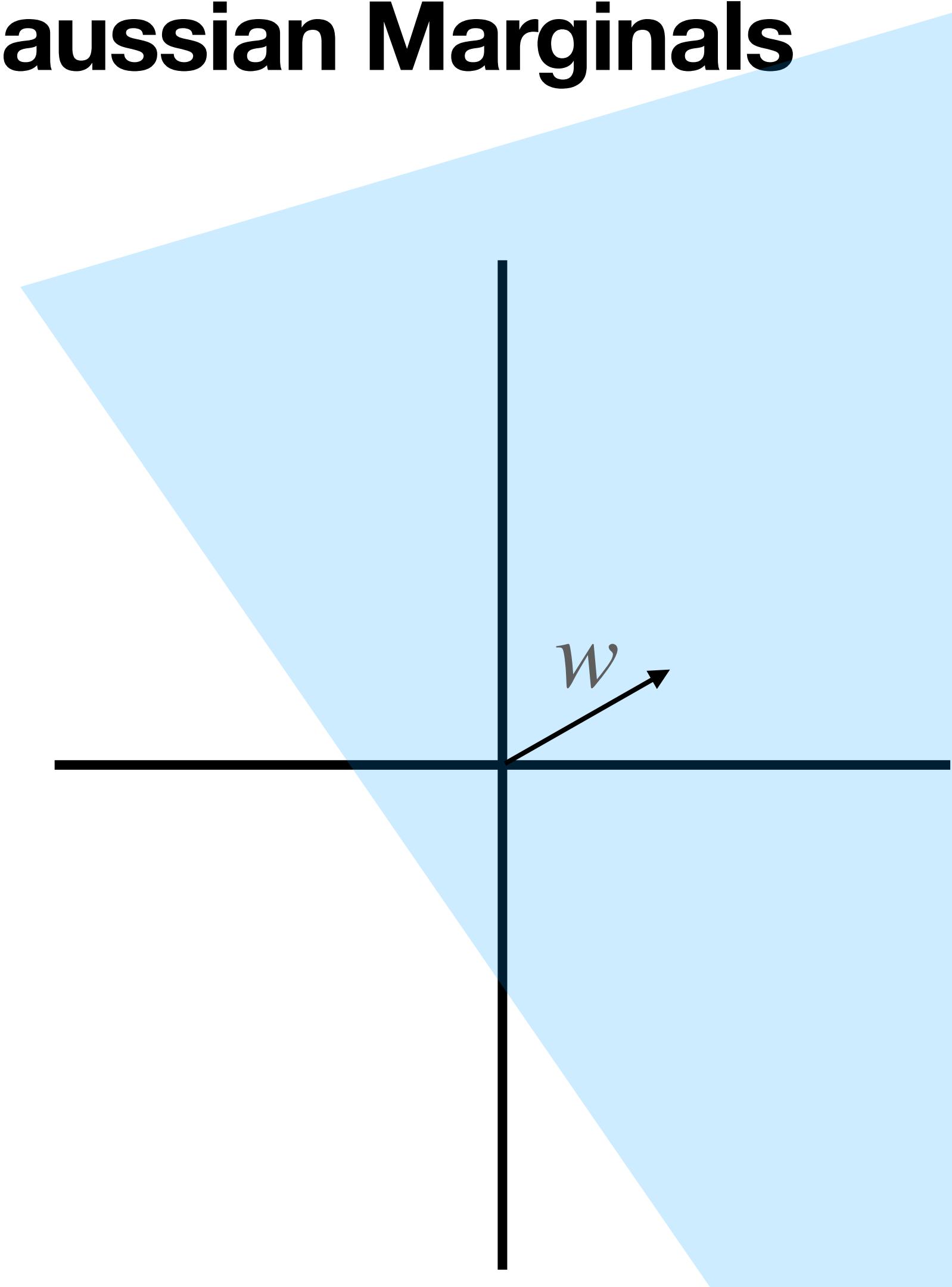
- **Problem:** Hard to prove computational lower bounds of agnostic learning.
- **Idea:** Use a restricted model for lower bounds: queries instead of samples.
- A **statistical query** takes (g, τ) with $g : \mathbb{R}^n \times \{\pm 1\} \rightarrow [-1, 1]$, tolerance $\tau > 0$ and returns $q \in \mathbb{E}[g(x, y)] \pm \tau$.
- Every SQ algorithm can be simulated by a sample-based algorithm:
 - For $m = \ln(2/\delta)/(2\tau^2)$ samples, $\frac{1}{m} \sum_{i=1}^m g(x_i, y_i) \in \mathbb{E}[g(x, y)] \pm \tau$ with probability $\geq 1 - \delta$.
- Many sample-based algorithms can be simulated with SQs: (e.g. gradient descent, polynomial regression)
 - But, not all! Parities are PAC-learnable, but not SQ-learnable.

Near-Optimal Statistical Query **Lower Bounds** for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals

- For hypothesis class \mathcal{H} , show that any agnostic SQ algorithm \mathcal{A} learning \mathcal{H} either
 1. makes $\geq \text{superpoly}(n)$ queries (); or
 2. makes at least one query of tolerance $\tau \leq 1/\text{superpoly}(n)$.

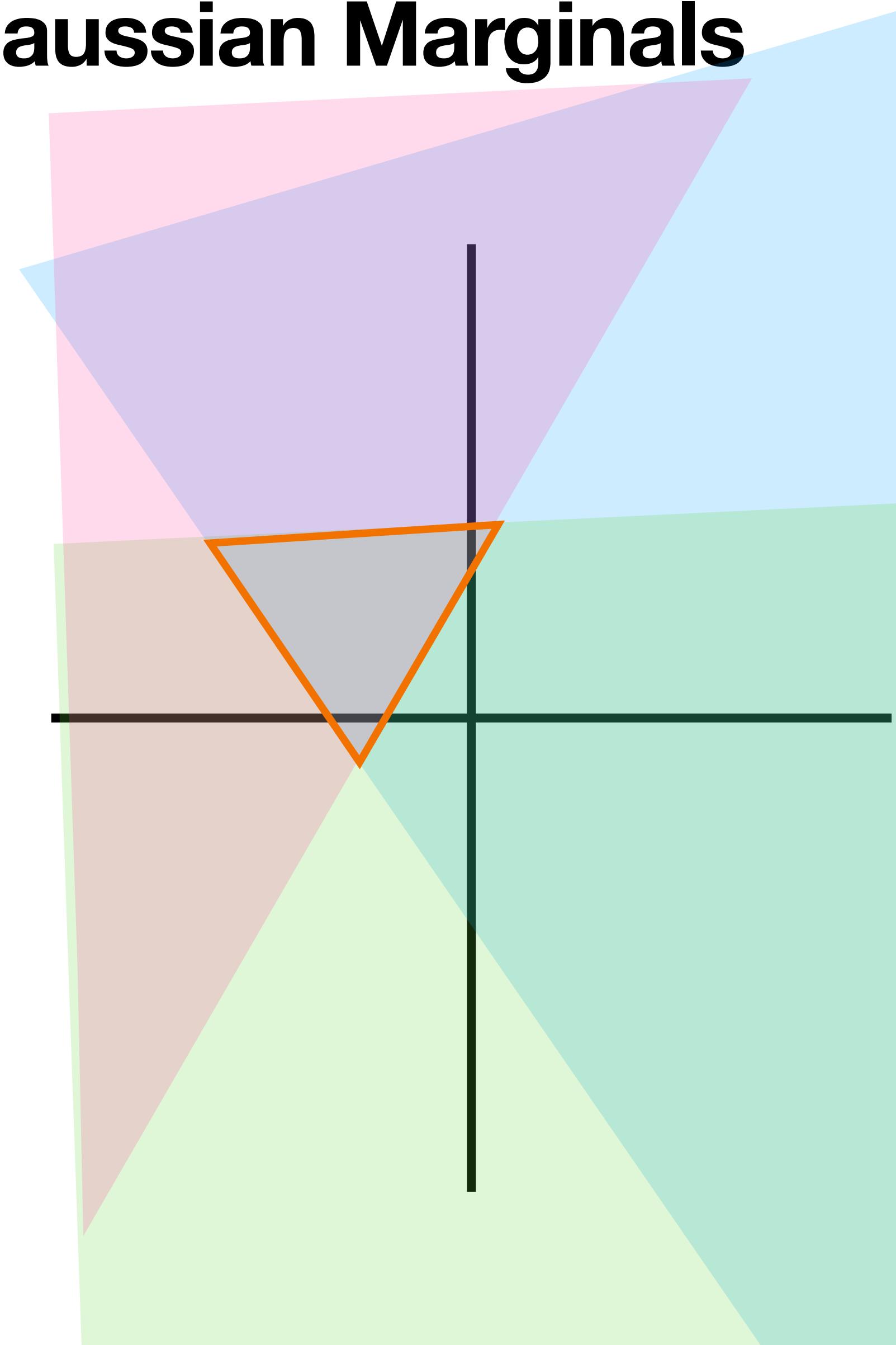
Near-Optimal Statistical Query Lower Bounds for Agnostically Learning **Intersections of Halfspaces** with Gaussian Marginals

- A **halfspace** is a function $h(x) = \text{sign}(w^T x - b)$.



Near-Optimal Statistical Query Lower Bounds for Agnostically Learning **Intersections of Halfspaces** with Gaussian Marginals

- A **halfspace** is a function $h(x) = \text{sign}(w^T x - b)$.
- \mathcal{H}_k is the class of all intersections of k halfspaces.
 - $f(x) = \min_{i \in [k]} \text{sign}(w_i^T x - b_i)$
 $= 2 \prod_{i=1} 1\{w_i^T x \geq b_i\} - 1$



Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces **with Gaussian Marginals**

- Features are drawn from a multivariate Gaussian distribution: $x \sim \mathcal{N}(0, I_n)$.

Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals

- [Klivans, O'Donnell, Servedio 2008] \mathcal{H}_k can be agnostically learned to accuracy ϵ with an L^1 polynomial approximation algorithm with $n^{O(\log k)}$ samples.
- Can be implemented as an SQ algorithm that makes $n^{O(\log k)}$ queries of tolerance $n^{-O(\log k)}$.
- **Question:** Is this dependence on k optimal?

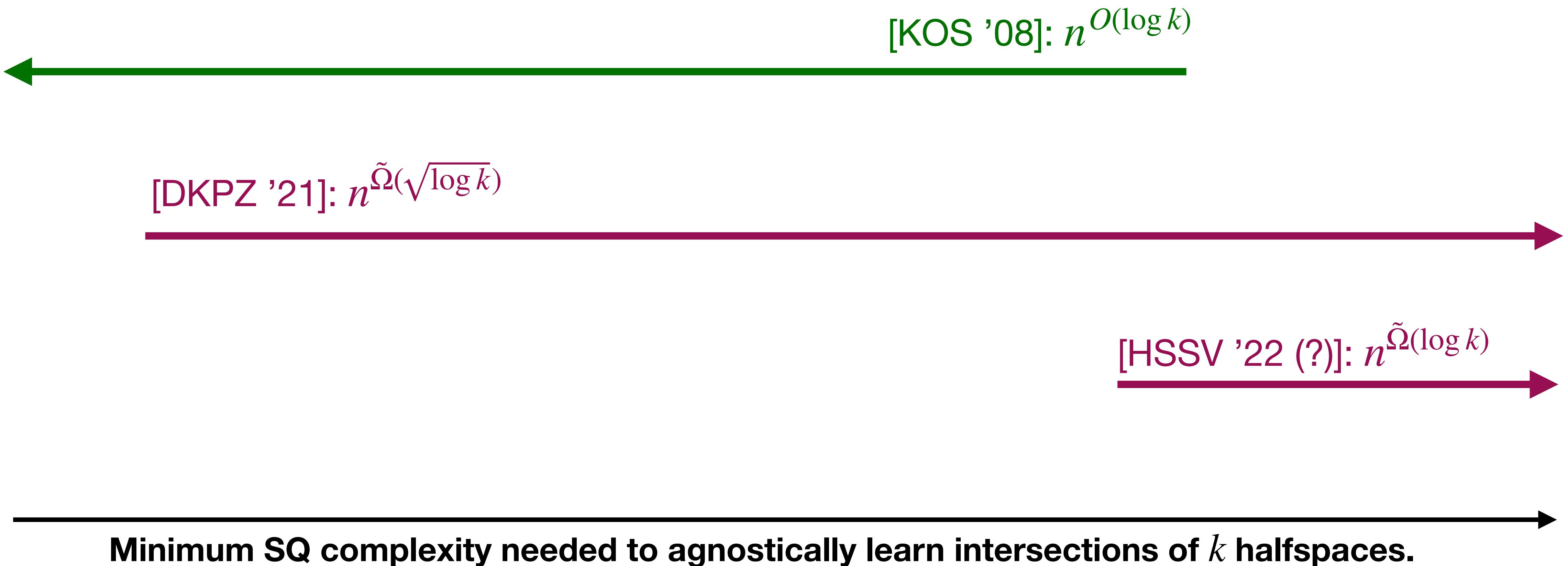
Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals (Ctd.)

- **[KOS 2008]** \mathcal{H}_k can be agnostically learned to accuracy ϵ with $n^{O(\log k)}$ queries of tolerance $n^{-O(\log k)}$.
- **[Diakonikolas, Kane, Pittas, Zarifis 2021]** To agnostically learn \mathcal{H}_k to accuracy ϵ , either $2^{n^{0.1}}$ queries are needed or at least one query of tolerance $\leq n^{-\tilde{\Omega}(\sqrt{\log k})}$ is necessary.
- **Question:** Which bound has the correct dependence on k ?

Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals (Ctd.)

- **[KOS 2008]** \mathcal{H}_k can be agnostically learned to accuracy ϵ with $n^{O(\log k)}$ queries of tolerance $n^{-O(\log k)}$.
- **[DKPZ 2021]** Requires either $2^{n^{0.1}}$ queries or at least one query of tolerance $\leq n^{-\tilde{\Omega}(\sqrt{\log k})}$.
- **Our results:** Requires either $2^{n^{0.1}}$ queries or at least one query of tolerance $\leq n^{-\tilde{\Omega}(\log k)}$.

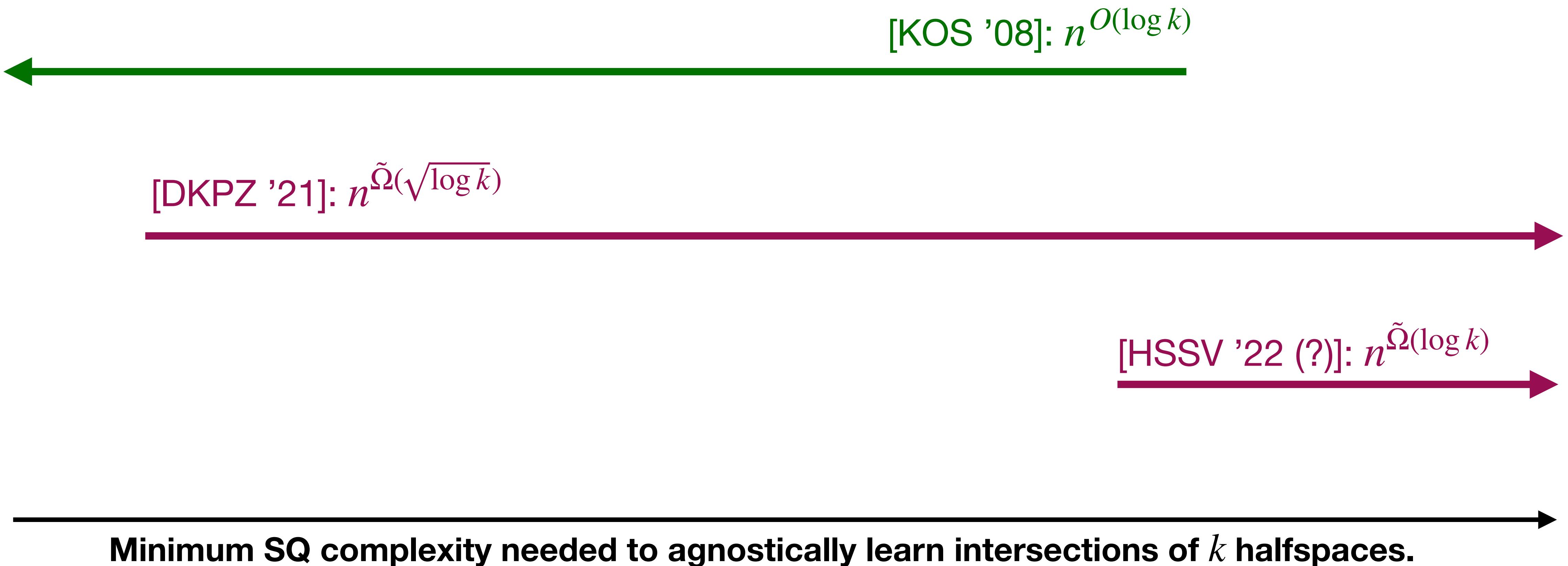
Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals (Ctd.)



Related problems

- Realizable learning \mathcal{H}_k with Gaussian marginals:
 - **[KOS '08]** $n^{O(\log k)}$ samples (L^1 polynomial approximation)
 - **[Vempala '10]** $\text{poly}(n, k) + k^{O(\log k)}$ (PCA approach)
- Agnostically learning halfspaces (\mathcal{H}_1) with Gaussian marginals:
 - **[KKMS '08]** $n^2 \log(n)/\epsilon^2$ samples
 - **[Ganzburg '02] + [DKPZ '21]** $n^{\Omega(1/\epsilon^2)}$ SQ complexity
- Learning \mathcal{H}_k in distribution-free setting thought to be hard:
 - **[Klivans, Sherstov '06]** Cryptographic hardness results
 - **[Sherstov '13]** No efficient algos with polynomial threshold hypotheses for \mathcal{H}_2

Near-Optimal Statistical Query Lower Bounds for Agnostically Learning Intersections of Halfspaces with Gaussian Marginals (Ctd.)



Proving an SQ lower bound

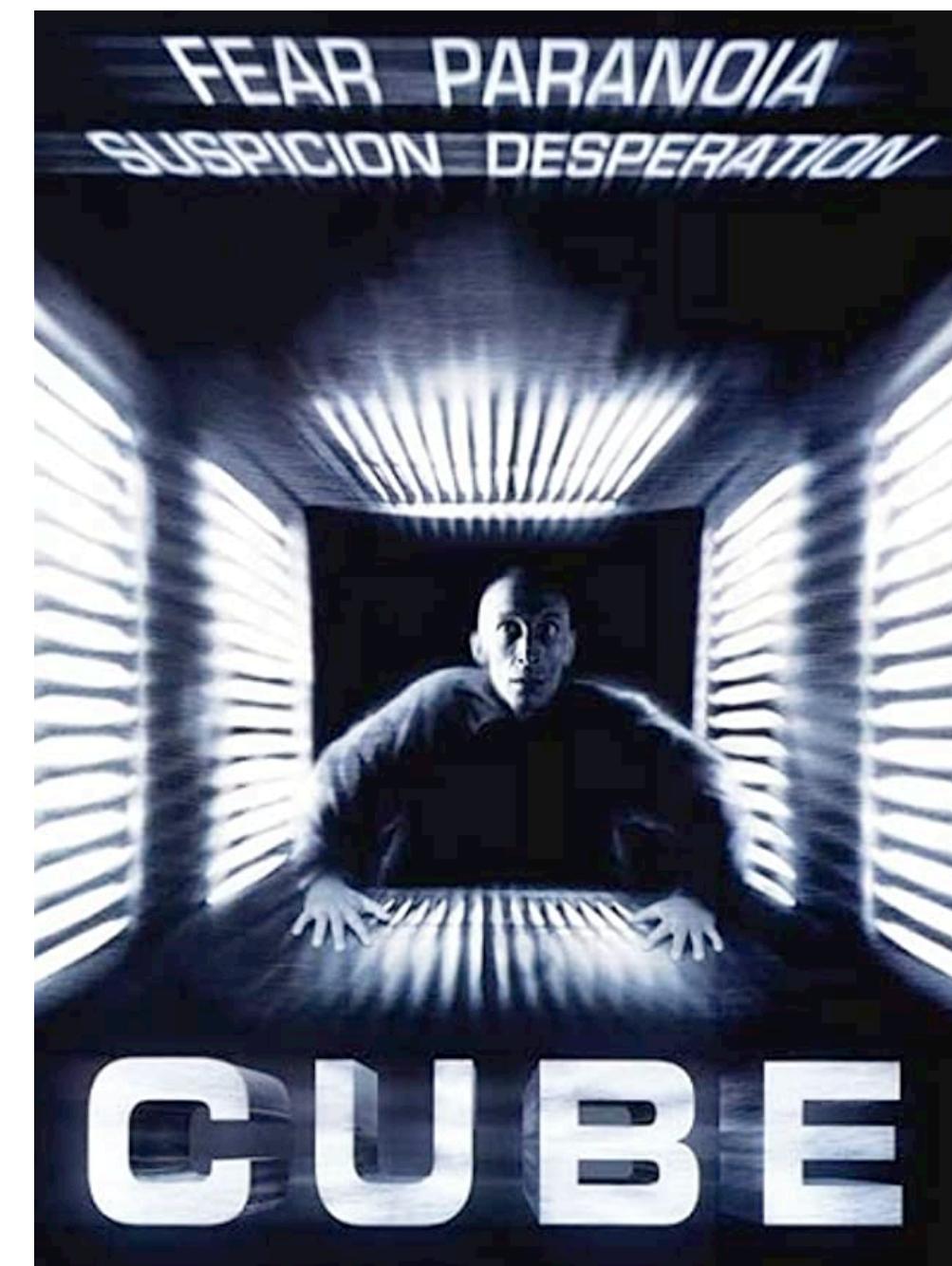
Reverse-engineering our result

- If there exists $h_1, \dots, h_m \in \mathcal{H}$ with $|\langle h_i, h_j \rangle| = |\mathbb{E}[h_i(x)h_j(x)]| \leq 1/m$, then the **SQ-dimension** of \mathcal{H} is at least m .
- Any SQ algorithm that learns \mathcal{H} to error $1/2 - m^{-1/3}$ must make either $\Omega(m^{1/3})$ queries or at least one query of tolerance $O(m^{-1/3})$.
- **Intuition:** If $f : \mathbb{R}^k \rightarrow \{-1, 1\}$ for $k \ll n$ cannot be approximated by low-degree polynomials, then $\mathcal{H} = \{x \mapsto f(Wx) : W \in \mathbb{R}^{k \times n}, WW^T = I\}$ has high SQ-dimension.
- Suffices to show existence of intersection of $\Theta(k)$ halfspaces f that is nearly orthogonal to low-degree polynomials.

SQ lower bounds on $\{\pm 1\}^n$

[Dachman-Soled, Feldman, Tan, Wan, Wimmer 2014]

- $f : \{\pm 1\}^k \rightarrow [-1,1]$ is **d -resilient** if for all $p \in \mathcal{P}_d$ (d -degree polys), $\langle f, p \rangle = 0$.
- f is **α -approximately d -resilient** if there exists $g : \{\pm 1\}^k \rightarrow [-1,1]$ such that $\|f - g\|_1 = \mathbb{E}_{x \sim \text{Unif}(\{\pm 1\}^k)} |f(x) - g(x)| \leq \alpha$ and g is d -resilient.
- **[DFTWW14, Thm 1.1]** For $k = n^{1/3}$, if $f : \{\pm 1\}^k \rightarrow \{\pm 1\}$ is α -approximately d -resilient, then agnostically learning $\mathcal{H} = \{f(x_S) : S \subset [n], |S| = k\}$ to excess error $(1 - \alpha)/2$ requires either $n^{\Omega(d)}$ queries or at least one query of tolerance $\leq n^{-\Omega(d)}$.
 - Lower bound on SQ dimension.
- **[DFTWW, Thm 1.6]** Tribes : $\{\pm 1\}^k \rightarrow \{\pm 1\}$ (read-once monotone DNF) is $O(k^{-1/3})$ -approximately $\Omega(\log(k)/\log \log k)$ -resilient.
 - Bounds on low-degree Fourier coefficients provide transform to resilient approximation.



Adaptation to $\mathcal{N}(0, I_n)$

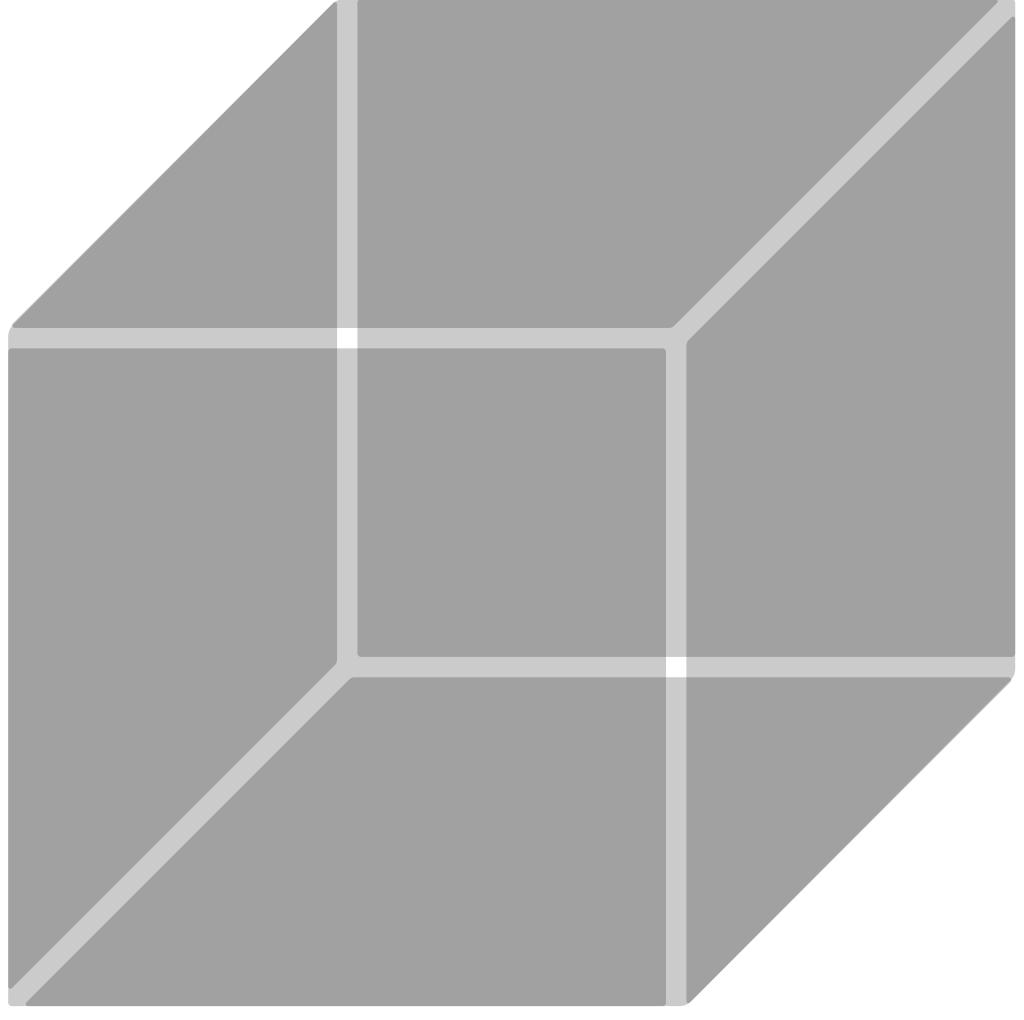
[Diakonikolas, Kane, Pittas, Zarifis 2021]

- $f : \mathbb{R}^k \rightarrow \{-1, 1\}$ is **α -approximately d -resilient** if there exists $g : \mathbb{R}^k \rightarrow [-1, 1]$ such that $\|f - g\|_1 \leq \alpha$ and $\text{Low}_d[g](x) = 0$.
- **[DKPZ '21, Prop 2.1]** f is α -approximately d -resilient iff $\|f - p\|_1 \geq 1 - \alpha$ for all $p \in \mathcal{P}_d$.
- **[DKPZ '21, Thm 1.4]** For $k = n^{0.1}$, if $\|f - p\|_1 \geq 1 - \alpha$ for all $p \in \mathcal{P}_d$, then learning $\mathcal{H} = \{x \mapsto f(Wx) : W \in \mathbb{R}^{k \times n}, WW^T = I\}$ to excess error $(1 - \alpha)/2$ requires either $2^{\Omega(n^{0.1})}$ queries or one query of tolerance $n^{-\Omega(d)}$.
- **[DKPZ '21, Thm 3.5]** For $d = \tilde{\Omega}(\sqrt{\log k})$, $\|f - p\|_1 \geq 0.1$ for all $p \in \mathcal{P}_d$ for some $f \in \mathcal{H}_k$.

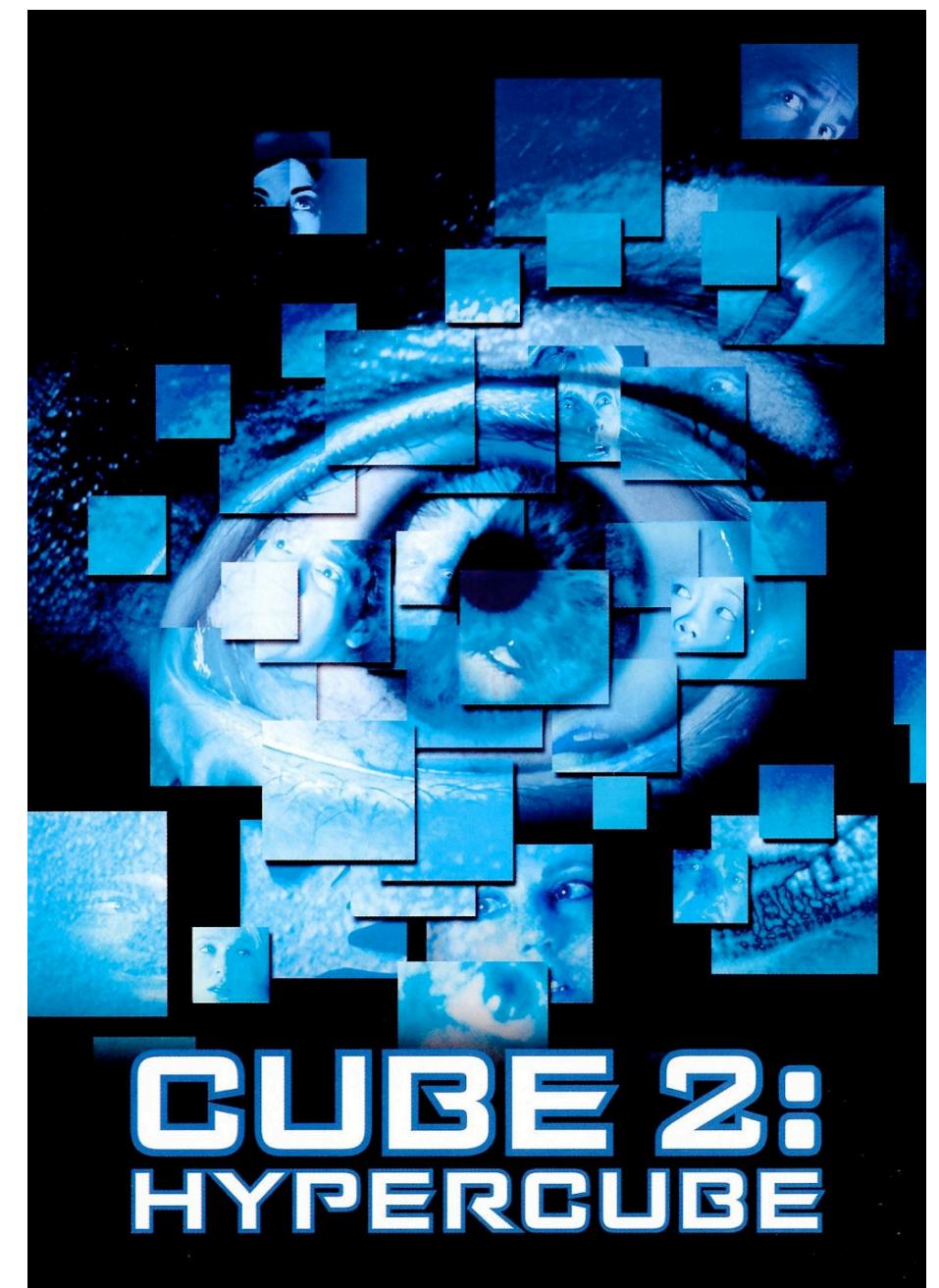
Our technical contributions:

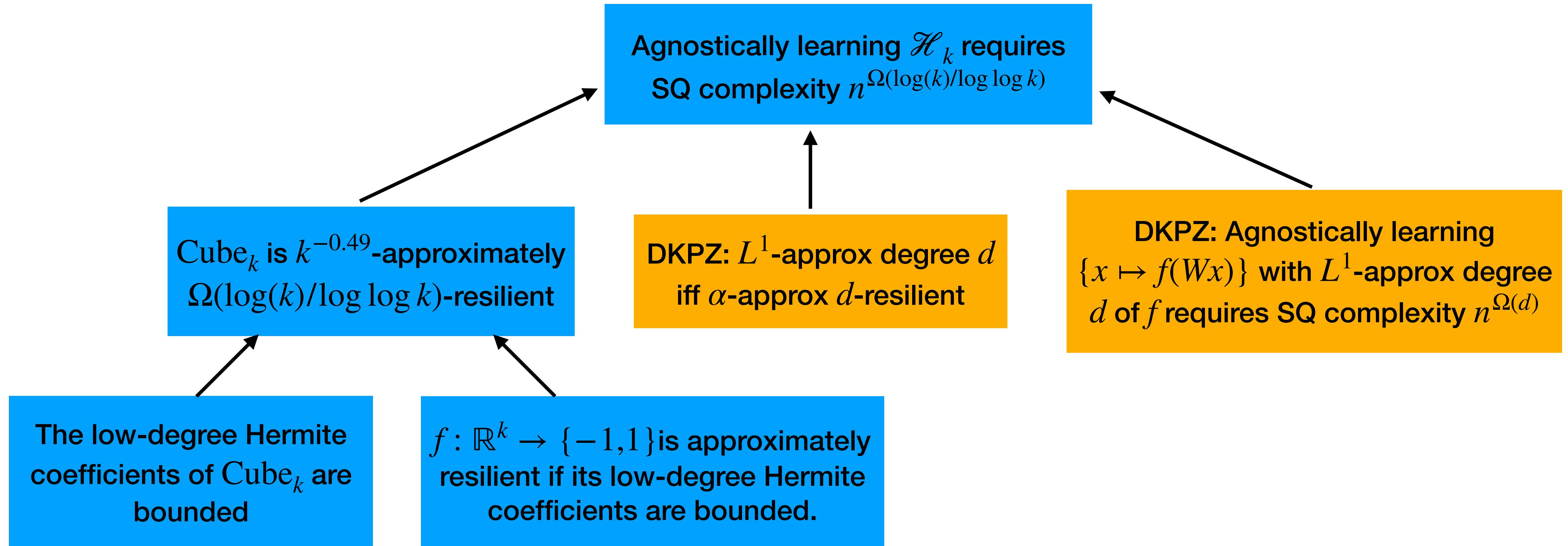
Improving the approximate resilience bound

- $\text{Cube}_k(x) = \text{sign}(\theta_k - \|x\|_\infty) = 2 \max_{i \in [k]} \mathbf{1}\{|x_i| \leq \theta_k\} - 1.$
 - $\theta_k = \Theta(\sqrt{\log k})$ chosen to have $\mathbb{E}_{x \sim \mathcal{N}(0, I_k)}[\text{Cube}_k(x)] = 0.$
- **Lemma:** Cube_k is $k^{-0.49}$ -approximately $\Omega(\log(k)/\log \log k)$ -resilient.
- **Theorem:** For $k = O(n^{0.49})$, agnostically learning $\mathcal{H} = \{x \mapsto \text{Cube}_k(Wx) : W \in \mathbb{R}^{k \times n}, WW^T = I\} \subset \mathcal{H}_{2k}$ to excess error $(1 - k^{-0.49})/2$ requires either $2^{O(n^{0.1})}$ queries or one query of tolerance $n^{-\Omega(\log(k)/\log \log k)}$.



$$2\theta_k$$





Hermite polynomials

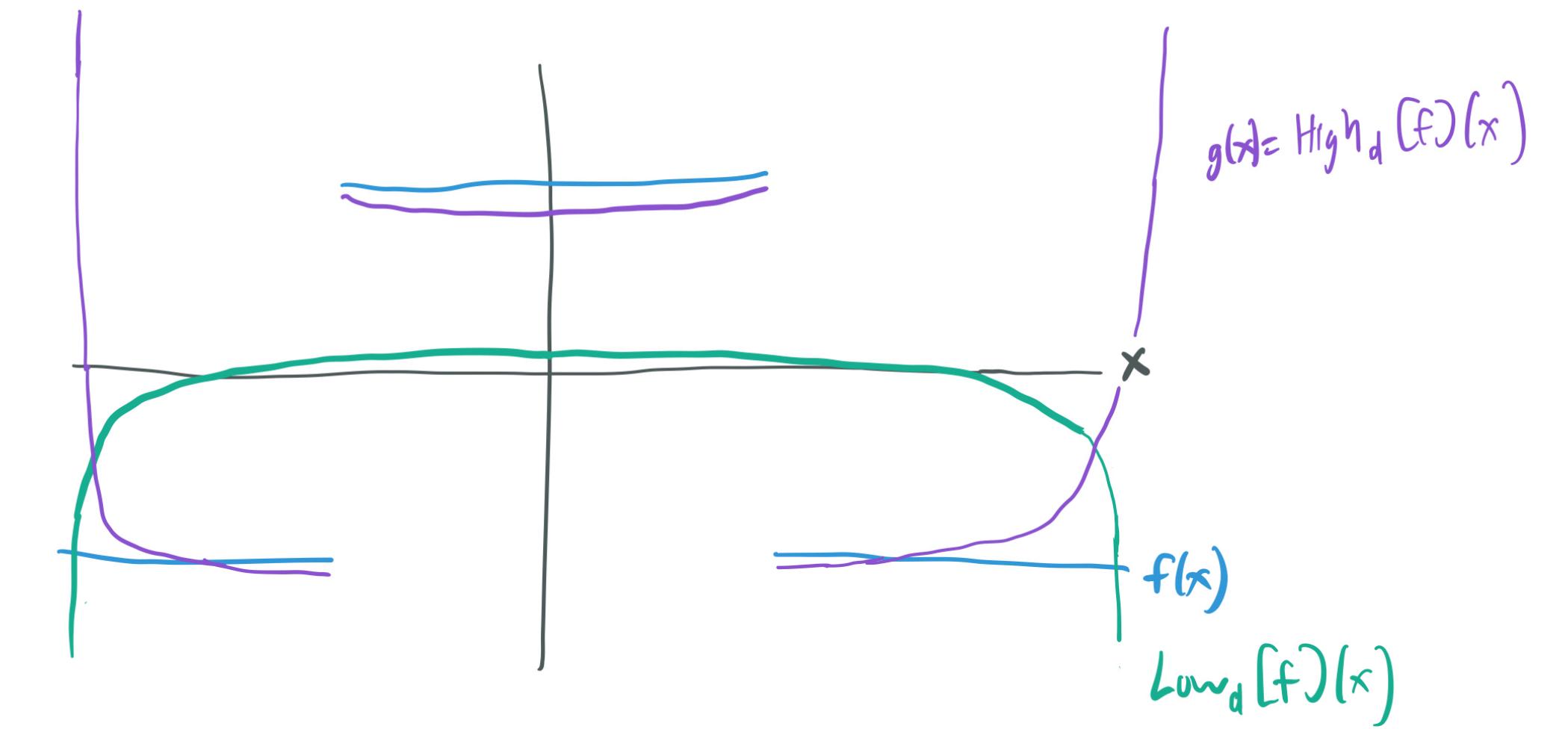
- Multivariate probabilist's Hermite polynomials $\{H_J\}_{J \in \mathbb{N}^k}$ is an orthogonal basis for $L^2(\mathcal{N}(0, I_k))$:
 - $\langle H_J, H_{J'} \rangle = \mathbb{E}_{x \sim \mathcal{N}(0, I_k)} [H_J(x) H_{J'}(x)] = J_1! \cdot \dots \cdot J_k! \mathbf{1}\{J = J'\}$
- Hermite representation: $f(x) = \sum_{J \in \mathcal{N}^k} \tilde{f}(J) H_J$, for $\tilde{f}(J) = \langle f, H_J \rangle / \sqrt{J!}$
- Decomposition of f :
 - $\text{Low}_d[f](x) = \sum_{|J| \leq d} \tilde{f}(J) H_J(x)$
 - $\text{High}_d[f](x) = f(x) - \text{Low}_d[f](x) = \sum_{|J| > d} \tilde{f}(J) H_J(x)$

Bound on low-degree Hermite coefficients of cube

- $\text{Cube}_k(x) = \text{sign}(\theta_k - \|x\|_\infty)$
- **Lemma:** $\|\text{Low}_d[\text{Cube}_k]\|_2^2 = \sum_{|J| \leq d} \widetilde{\text{Cube}_k}(J)^2 \leq \frac{(4 \ln k)^d}{k}$
 - For $d = \ln(k)/400 \ln \ln k$, gives $\|\text{Low}_d[\text{Cube}_k]\|_2^2 \leq k^{-0.99}$
 - Proof by exact computation of univariate Hermite coefficients of $t \mapsto \mathbf{1}\{|t| \leq \theta_k\}$, bound by Stirling inequality and $\theta_k \in [\sqrt{2 \ln k - \ln(2 \ln k)}, \sqrt{2 \ln k}]$

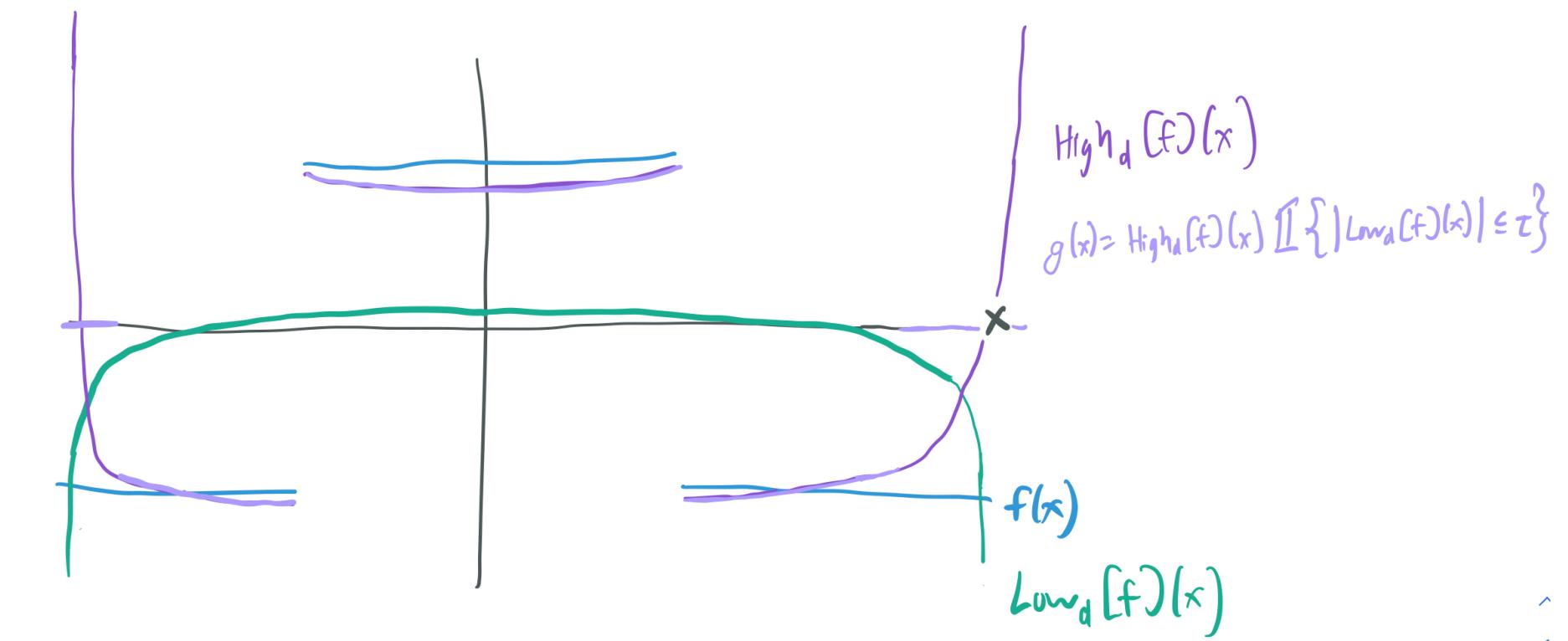
Approximate Resilience: Truncation Transform

- $f : \mathbb{R}^k \rightarrow \{-1, 1\}$ is **α -approximately d -resilient** if there exists $g : \mathbb{R}^k \rightarrow [-1, 1]$ such that $\|f - g\|_1 \leq \alpha$ and $\text{Low}_d[g](x) = 0$.
- **Goal:** Transform f with small $\|\text{Low}_d[f]\|_2$ to obtain **bounded** g that **approximates** f and is **uncorrelated** to low-degree polynomials.
- **Idea #1:** $g(x) := \text{High}_d[f](x)$.
 - $\|g - f\|_1 \leq \|\text{Low}_d[f]\|_2$ 😊
 - $\text{Low}_d[g] = 0$ 😎
 - g is not bounded 💀



Approximate Resilience: Truncation Transform

- $f : \mathbb{R}^k \rightarrow \{-1, 1\}$ is **α -approximately d -resilient** if there exists $g : \mathbb{R}^k \rightarrow [-1, 1]$ such that $\|f - g\|_1 \leq \alpha$ and $\text{Low}_d[g](x) = 0$.
- **Goal:** Transform f to obtain **bounded** g that **approximates** f and is **uncorrelated** to low-degree polynomials.
- **Idea #2:** $g = \text{High}_d[f](x) \cdot \mathbf{1}\{|\text{Low}_d[f](x)| \leq \eta\}$
 - For large η , $\|g - f\|_1 \leq 2\|\text{Low}_d[f]\|_2$ 😊
 - For large η , $\|\text{Low}_d[g]\|_2 \leq \|\text{Low}_d[f]\|_2/a$ 🙄
 - $\|g\|_\infty \leq 1 + \eta$ 😞

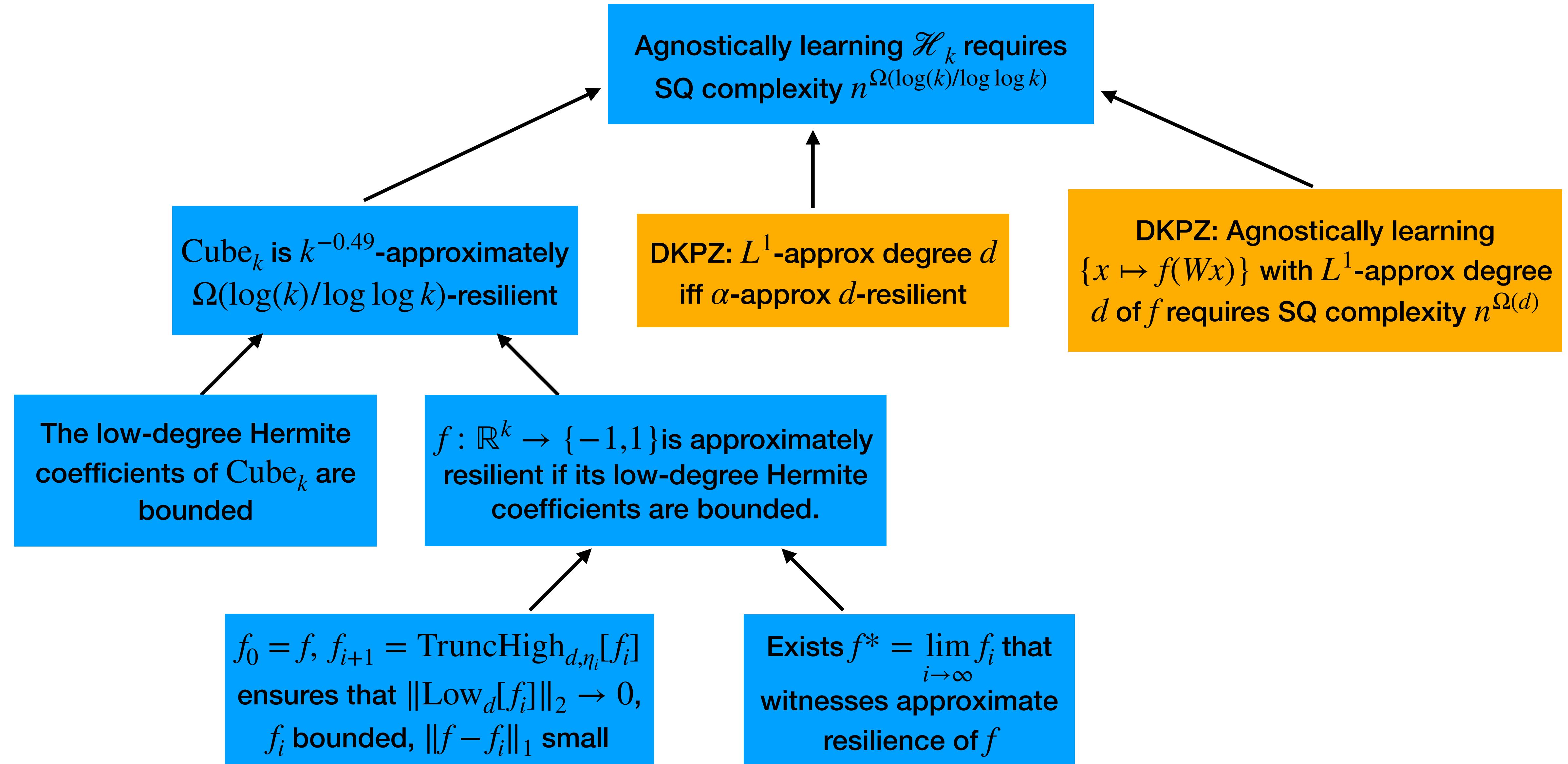


Approximate Resilience: Truncation Transform

- $f : \mathbb{R}^k \rightarrow \{-1, 1\}$ is **α -approximately d -resilient** if there exists $g : \mathbb{R}^k \rightarrow [-1, 1]$ such that $\|f - g\|_1 \leq \alpha$ and $\text{Low}_d[g](x) = 0$.
- **Goal:** Transform f to obtain **bounded** g that **approximates** f and is **uncorrelated** to low-degree polynomials.
- **[DFTWW '14]**
 - $h(x) = \text{High}_d[\text{High}_d[f] \cdot \mathbf{1}\{|\text{Low}_d[f]| \leq \eta\}]$
 - $g(x) = h(x)/\|h\|_\infty$.
 - 

Approximate Resilience: Truncation Transform

- $f : \mathbb{R}^k \rightarrow \{\pm 1\}$ is **α -approximately d -resilient** if there exists $g : \mathbb{R}^k \rightarrow [-1,1]$ such that $\|f - g\|_1 \leq \alpha$ and $\text{Low}_d[g](x) = 0$.
- $\text{TruncHigh}_{d,\eta}[f] = \text{High}_d[f](x) \cdot \mathbf{1}\{|\text{Low}_d[f](x)| \leq \eta\}$
- Let $f_0 := f$ and $f_{i+1} = \text{TruncHigh}_{d,\eta_i}[f_i]$ for $i \rightarrow \infty$.
- For some decaying η_i and $\alpha = k^{0.49}$, have (1) $\|f_{i+1}\|_\infty \leq \|f_i\|_\infty + \alpha/(3 \cdot 2^{i+1})$,
(2) $\lim_{i \rightarrow \infty} \|\text{Low}_d[f_i]\| = 0$, and (3) $\|f_{i+1} - f_i\|_1 \leq \alpha/(3 \cdot 4^i)$.
- By limit argument, exists f^* with (1) $\|f^*\|_\infty \leq 1 + \alpha/3$, (2) $\text{Low}_d[f^*] = 0$,
and (3) $\|f - f^*\|_1 \leq 2\alpha/3$. Let $g := f^*/\|f^*\|_\infty$.



What else is there?

- Second proof for larger $k = 2^{O(n^{0.245})}$ (rather than $k = O(n^{0.49})$) bounds on L^1 approximate degree of random intersection of halfspaces.
 - Based on hardness of weak-learning intersections of halfspaces with membership queries **[De, Servedio 2021]**
- Dependence on accuracy ϵ : $n^{\Omega(\log(k)/\log \log k + 1/\epsilon^2)}$ bound by augmenting construction with a single centered halfspace **[Ganzburg 2002]**
- Optimality of learning families with *Gaussian surface area* $\leq s$ with L^1 polynomial approximation: SQ complexity of $n^{\Omega(s^2/\log s)}$, vs $n^{O(s^2)}$ **[KOS 2008]**

Thanks!