

11/8/23

# Lecture 18: Statistical Query Model (cont'd)

SQ dimension (recipe for supervised problems)

Def: Class of functions  $\mathcal{F} = \{f: \mathbb{R} \rightarrow [-1, 1]\}$  has

SQ dimension  $\geq D$  w.r.t  $q$  if  $\exists f_1, \dots, f_D \in \mathcal{F}$   
s.t.  $\forall i \neq j$

$$\left| \mathbb{E}_{x \sim q} [f_i(x) f_j(x)] \right| \leq \frac{1}{D}$$

Thm: If  $\mathcal{F}$  has SQ dimension  $D$ , then CSQ learning requires tolerance  $\tau$  or  $\Omega(D\tau^2)$  queries

Pf: For convenience, define  $\langle f, g \rangle \triangleq \mathbb{E}[f(x)g(x)]$ .

wlog let  $\mathcal{F} = \{f_1, \dots, f_D\}$ .

Given query  $\phi: \mathbb{R}^d \rightarrow [-1, 1]$ , let

$$A^+ \triangleq \left\{ f \in \mathcal{F} : \langle f, \phi \rangle \geq \tau \right\}$$

By Cauchy-Schwartz:

$$\begin{aligned} \left\langle \phi, \sum_{f \in A^+} f \right\rangle^2 &\leq \underbrace{\|\phi\|^2}_{\leq 1} \cdot \left\| \sum_{f \in A^+} f \right\|^2 \\ &\leq \sum_{f \in A^+} \|f\|^2 + \frac{|A^+|(|A^+| - 1)}{D} \\ &\leq \frac{|A^+|^2}{D} + |A^+| \end{aligned}$$

But  $\langle \phi, \sum_{f \in A^+} f \rangle \geq \tau |A^+|$  by defn., so

$$\tau^2 |A^+|^2 \leq \frac{|A^+|^2}{D} + |A^+|$$

$$\Rightarrow |A^+| \leq \frac{D}{D\tau^2 - 1} \leq O\left(\frac{1}{\tau^2}\right)$$

Similarly, for  $A^- \triangleq \{f \in \mathcal{F} : \langle f, \phi \rangle \leq -\tau\}$ ,  
 $|A^-| \leq O\left(\frac{1}{\tau^2}\right)$ .

So regardless of  $\phi$ , all but  $O\left(\frac{1}{\tau^2}\right)$  many  
 $f$ 's consistent w/ the answer 0, so  
need  $\Omega(D\tau^2)$  queries.  $\square$

---

SQ dimension bound for 1-hidden-layer MLPs:

[Goel - Gollakota - Jin - Karmalkar - Klivans '20]:

Let  $S \subseteq [d]$  be of size  $m = \lg_2 k$ .

Given  $w \in \{\pm 1\}^m$ , define  $w^{(S)} \in \mathbb{J}^{d-1}$  by

$$w_i^{(S)} = \begin{cases} w_i / \sqrt{m} & \text{if } i \in S \\ 0 & \text{o.w.} \end{cases}$$

Define  $f_S : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$f_S(x) \triangleq \sum_{w \in \{\pm 1\}^m} \left( \prod_{i=1}^m w_i \right) \cdot \text{relu}(\langle w^{(S)}, x \rangle)$$

(Note: this is one-hidden-layer MLP w/ size  $2^m = k$ )

Claim: If  $q$  is sign-symmetric dist. over  $\mathbb{R}$

(ie. for any  $x$ ,  $x$  and  $-x$  are equally likely under  $q$ ), then  $\text{span dim of } \{f_S\}$

is  $\geq \binom{d}{m}$ .

Pf: For any distinct  $S, T$  of size  $m$ ,

$$\text{w.t.s } \langle f_S, f_T \rangle_q = 0$$

For any  $z \in \{\pm 1\}^d$ , note that

$$\begin{aligned}
 f_S(x \odot z) &= \sum_w \prod_i w_i \cdot \underbrace{\text{relu}(\langle w^{(S)}, x \odot z \rangle)}_{\parallel \text{relu}(\langle w^{(S)} \odot z, x \rangle)} \\
 &= \sum_{w'} \prod_i w'_i \cdot z_S \cdot \text{relu}(\langle w'^{(S)}, x \rangle) \\
 &= z_S \cdot f_S(x).
 \end{aligned}$$

by sign symmetry

$$\begin{aligned}
 \text{so } \langle f_S, f_T \rangle_q &\stackrel{\downarrow}{=} \int_x \int_z [f_S(x \odot z) f_T(x \odot z)] \\
 &= \int_x \int_z [f_S(x) f_T(x) z_S z_T] \\
 &= \int_x [f_S(x) f_T(x)] \cdot \underbrace{\int_z [z_S z_T]}_0. \quad \square \\
 &\quad \text{if } S \neq T
 \end{aligned}$$

NB: Technically, also have to make sure  $\{f_S\}$  are nonzero functions! See paper for this calculation (Hermite analysis).

(Diakonikolas - Kane - Kontonis - Zarifis '20)

Showed stronger lbd ( $d^{-\Omega(k)}$  queries or  $2^{\text{poly}(d)}$  tolerance)

using different instance of approximately orthogonal

functions: given 2D subspace  $U \subseteq \mathbb{R}^d$ ,

$$f_U(x) = h(\langle u_1, x \rangle, \langle u_2, x \rangle)$$

where  $h(a, b) = \sum_{i=1}^k (-1)^i \left\langle \left( \cos\left(\frac{2\pi}{k}\right), \sin\left(\frac{2\pi}{k}\right) \right), (a, b) \right\rangle$ .

and  $u_1, u_2$  are basis for  $U$ .

idea:  $\deg \leq \frac{k+1}{2}$  Hermite coeffs of  $f_U$  are 0, so

$$\langle f_U, f_V \rangle_{N(0, I_d)} = \left\langle f_U^{\frac{k+1}{2}}, f_V^{\frac{k+1}{2}} \right\rangle$$

$$\leq \left( \text{"Correlation}(U, V) \right)^{\binom{n}{k}},$$

So if we take a packing of the space of 2D subspaces, we get the desired lbd on the SQ dimension.

---

### Statistical dimension :

Def: Set of distributions  $\mathcal{T}$  has statistical dimension  $\geq \Delta$  w.r.t. reference distribution  $D$  and avg. correlation  $\gamma$  if  $\forall \mathcal{T}$  of size  $\geq |\mathcal{T}|^2 / \Delta$ ,

$$\rho_D(\mathcal{T}) \stackrel{\Delta}{=} \frac{1}{|\mathcal{T}|^2} \sum_{D_i, D_j \in \mathcal{T}} \langle D_i, D_j \rangle_D \leq \gamma,$$

where  $\langle D_i, D_j \rangle_D = \mathbb{E}_{x \sim D} \left[ \left( \frac{D_i(x)}{D(x)} - 1 \right) \left( \frac{D_j(x)}{D(x)} - 1 \right) \right]$

Thm: If  $\mathcal{T}^a$  has statistical dimension  $\geq \Delta$ , then learning dist's in  $\mathcal{T}^a$  w/ SQ alg. requires either  $\Omega(\Delta)$  queries or  $\sqrt{\gamma}$  tolerance.

Pf: Given query  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$ , define

$$A^+ \stackrel{\circ}{=} \left\{ D' \in T^0 \text{ s.t. } \mathbb{E}_{D'}[\phi] \geq \mathbb{E}_D[\phi] + \sqrt{\gamma} \right\}$$

$$A^- \stackrel{\circ}{=} \left\{ D' \in T^0 \text{ s.t. } \mathbb{E}_{D'}[\phi] \leq \mathbb{E}_D[\phi] - \sqrt{\gamma} \right\},$$

i.e. set of dist's in  $T^0$  s.t. answering  $\phi$  w/  $\mathbb{E}_D[\phi]$  is incorrect. w.t.s  $|A^+|, |A^-|$  small.

$$\begin{aligned} \mathbb{E}_{x \sim D} \left[ \phi(x) \cdot \sum_{D' \in A^+} \left\{ \frac{D'(x)}{D(x)} - 1 \right\} \right]^2 &\leq \underbrace{\mathbb{E}_x \left[ \phi(x)^2 \right]}_{\leq 1} \cdot \underbrace{\mathbb{E}_x \left[ \left( \sum_{D' \in A^+} \frac{D'(x)}{D(x)} - 1 \right)^2 \right]}_{=} \\ &= \sum_{D', D'' \in A^+} \langle D', D'' \rangle_D \\ &\leq |A^+|^2 \cdot \rho_D(A^+) \end{aligned}$$

Remains to lower bound

$$\mathbb{E}_{x \sim D} \left[ \phi(x) \cdot \sum_{D' \in A^+} \left\{ \frac{D'(x)}{D(x)} - 1 \right\} \right]^2 \quad (\dagger)$$

$$\text{Note: } \mathbb{E}_{x \sim D} \left[ \phi(x) \cdot \left\{ \frac{D'(x)}{D(x)} - 1 \right\} \right] = \mathbb{E}_{x \sim D'}[\phi(x)] - \mathbb{E}_{x \sim D}[\phi(x)]$$

$$\geq \sqrt{\gamma},$$

So  $(\delta) \geq |A^+|^2 \cdot \gamma$ , so

$$\int_D(A^+) \geq \gamma,$$

implying  $|A^+|/|T^0| \leq 1/\Delta$ .

Similarly,  $|A^-|/|T^0| \leq 1/\Delta$ .

So need  $\Omega(\Delta)$  queries b/c each one only rules out a  $1/\Delta$  fraction.  $\square$

---

## (Bonus material)

Moment-matching and SQ:

Consider  $P_v$  and  $P_{v'}$  (defined in slides)

for moment-matching distribution  $A$ . WLOG

Suppose  $v = (1, 0, \dots, 0)$ ,  $v' = (\cos\theta, \sin\theta, 0, \dots, 0)$ .

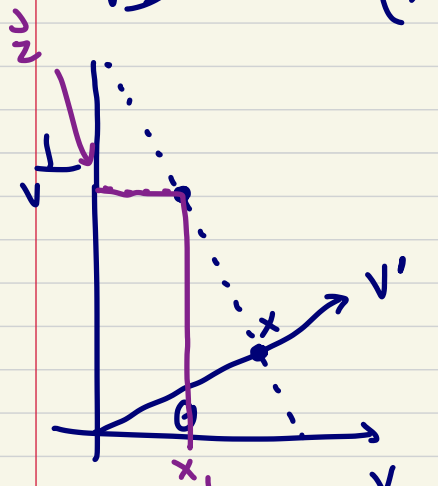
Claim: Let  $Q$  be dist over  $\mathbb{R}^d$

given by taking  $X \sim P_v$  and outputting  $v' \cdot X$ .



Then  $\langle Q, Q \rangle_{N(0,1)} \leq (\cos \theta)^{2(n+1)} \langle A, A \rangle_{N(0,1)}$

Pf:  $Q(x) = \int P_v(\vec{x}) d\vec{x}$



$\vec{x} \cdot (\vec{x}, v') = x$

Gaussian  $\downarrow$

$= \int A(x_1) \cdot G(\vec{z}) dx_1 d\vec{z}$

$(x_1, \vec{z})$ :

$\langle (x_1, \vec{z}), v' \rangle = x$

$\Downarrow$

$x_1 \cos \theta + z \sin \theta = x$

Scaled convolution of A and Gaussian

$= \int A(x_1) G(z_1) dx_1 dz_1$

$(x_1, z_1)$ :

$x_1 \cos \theta + z_1 \sin \theta = x$

$= \int_{g \sim N(0,1)} [A(\cos \theta x + \sin \theta g)]$

$\stackrel{D}{=} \int_{\theta} A(x)$

"Ornstein - noise Uhlenbeck operator"

Write  $\frac{A(x)}{G(x)} = \sum_{i=0}^{\infty} a_i \underbrace{He_i(x)}_{\text{probabilist's Hermite polynomial}} \cdot \frac{1}{\sqrt{i!}}$

note:  $\forall k \leq m$   
 $a_i = \int_{x \sim N(0,1)} \frac{A(x)}{G(x)} He_i(x) \cdot \frac{1}{\sqrt{i!}}$

(66)  $\stackrel{\text{b/c moment matching}}{=} \int A(x) He_i(x) \cdot \frac{1}{\sqrt{i!}} dx$   
 $\stackrel{\text{orthogonality}}{=} \int G(x) He_i(x) \cdot \frac{1}{\sqrt{i!}} dx \stackrel{!}{=} 0$

Mehler's lemma:  $U_{\theta}(He_i \cdot G)(x) =$

$\cos^i(\theta) He_i(x) G(x),$

i.e. (Hermite x Gaussian density) is eigenfunction of noise operator

So  $U_{\theta} A(x) = \sum_{i=0}^{\infty} a_i \cos^i(\theta) He_i(x) G(x) \frac{1}{\sqrt{i!}}$   
 $\stackrel{\text{(by (66))}}{=} \sum_{i=m+1}^{\infty} a_i \cos^i(\theta) He_i(x) G(x) \frac{1}{\sqrt{i!}}$

Intuition: "Smears out" noise high-degree info

$$\langle Q, Q \rangle_{N(0,1)} = \langle U_\theta A, U_\theta A \rangle_{N(0,1)}$$

$$= \sum_{i=m+1}^{\infty} a_i^2 \cos^{2i}(\theta)$$

$$\leq \cos^{2m+2}(\theta) \cdot \underbrace{\sum_{i=m+1}^{\infty} a_i^2}_{= \langle A, A \rangle_{N(0,1)}} \quad \square$$

We conclude that

$$\langle P_V, P_{V'} \rangle_{N(0, Id)} = \langle A, U_\theta A \rangle_{N(0,1)}$$

by the above

$$\leq \langle A, A \rangle_{N(0,1)}^{1/2} \cdot \langle U_\theta A, U_\theta A \rangle_{N(0,1)}^{1/2}$$

$$\leq \underbrace{\cos^{m+1}(\theta)}_{\text{exp. decaying in } m} \langle A, A \rangle_{N(0,1)} \quad \square$$