

sketching \rightarrow randomized LA

nonparametric estimation \rightarrow matrix \Rightarrow function

variance-reduced \rightarrow specific choice of sketching

1. sketching

1.1 matrix sketching

suppose matrix $B \in \mathbb{R}^{n_1 \times n_2}$ low-rank

① construct sketch $B' = BS$, $S \in \mathbb{R}^{n_2 \times k}$

$k \ll n_1, n_2$
sketching matrix.

$$\text{Range}(BS) = \text{Range}(B)$$

② orthogonalize

$$[Q, R] = \text{qr}(B') \quad / \quad [U, \Sigma, V] = \text{svd}(B)$$

\downarrow $\mathbb{R}^{n_1 \times k}$ \downarrow $\mathbb{R}^{n_1 \times k}$

$$\tilde{B} = QQ^T B$$

\downarrow
rank-k

$$\tilde{B} = UU^T B$$

\downarrow
rank-k

$$\|\tilde{B} - B\| \rightarrow \text{small}$$

$$P = QQ^T / UU^T$$

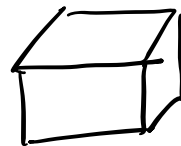
projection

How to get P

for continuous function?

1.2 tensor sketching.

$$B \in \mathbb{R}_{\text{index}(i_1, \dots, i_d)}^{n_1 \times n_2 \times \dots \times n_d} \rightarrow \text{low-rank.}$$



Step 1: $B^{(1)} = B[i_1; i_2, i_3, \dots, i_d] \in \mathbb{R}^{n_1 \times \prod_{j=2}^d n_j}$
 ↪ reshape d-1

use matrix sketching for $B^{(1)}$

$$S_1 \in \mathbb{R}^{\prod_{j=2}^d n_j \times k_1}$$

low-rank $Q_1 / U_1 \in \mathbb{R}^{n_1 \times k_1}$

$$P_1 = Q_1 Q_1^T / P_1 = U_1 U_1^T \rightarrow \text{projection}$$

Step 2: repeat the procedure, $\forall j=1, \dots, d$

$$B^{(j)} = B[i_j; i_1, i_2, \dots, i_{j-1}, i_{j+1}, \dots, i_d]$$

permuto

$$\stackrel{\text{reshape}}{=} B[i_j; i_1, i_2, \dots, i_{j-1}, i_{j+1}, \dots, i_d]$$

$$\in \mathbb{R}^{n_j \times \prod_{l \neq j} n_l}$$

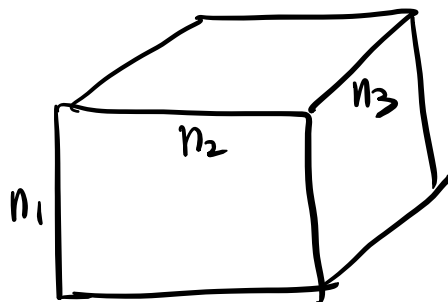
use matrix sketching for $B^{(j)}$ with S_j

get Q_j / U_j , projection $P_j = Q_j Q_j^T / U_j U_j^T$

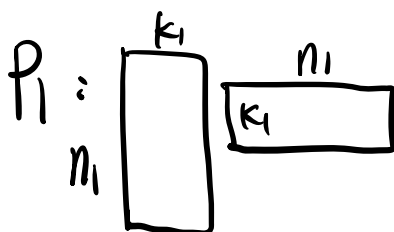
Step 3: $\tilde{B} = B \times_1 P_1 \times_2 P_2 \times \dots \times_d P_d$

$(B \times_1 P_1)(i_1, \dots, i_d) = \sum_j B(j, i_2, \dots, i_d) P_1(j, i_1)$

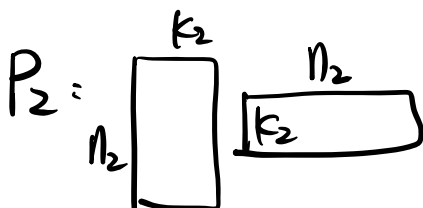
$B \in \mathbb{R}^{n_1 \times n_2 \times n_3}$



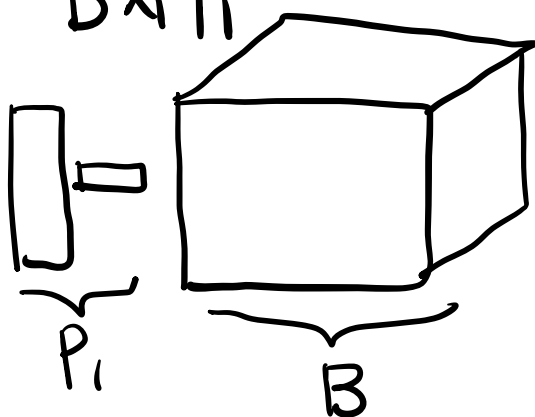
$P_1 \in \mathbb{R}^{n_1 \times n_1}$
rank- k_1



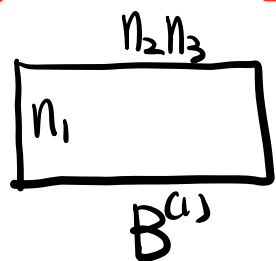
$P_2 \in \mathbb{R}^{n_2 \times n_2}$



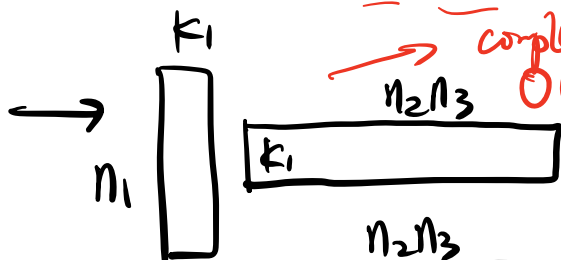
$B \times_1 P_1$



P_3 : complexity $O(k_1 n_2 n_3)$

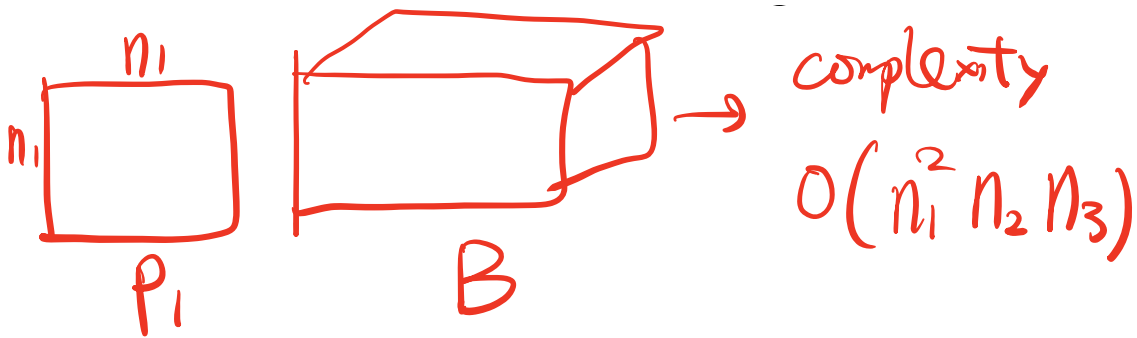


complexity $O(n_1 k_1 n_2 n_3)$

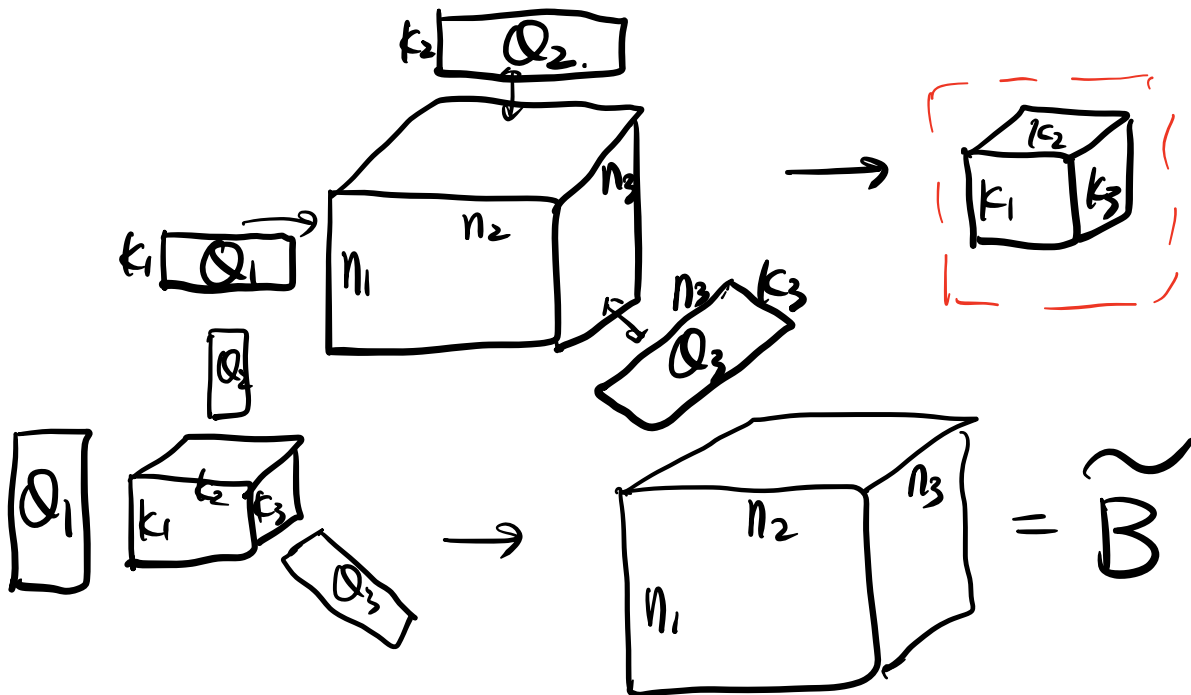


$k_1 \ll n_1$





$$B \times_1 P_1 \times_2 P_2 \times_3 P_3$$



$B \in \mathbb{R}^{n \times n \times \dots \times n}$ $O(n^d)$ $k \ll n$
 low-rank approximation. $O(dnk + k^d)$

2. Nonparametric estimation

2.1 Function representation

Basis

$f(x)$, $x \in \mathbb{R}^1$ $\{\phi_i(x)\}$ orthonormal basis function

Result:

$$f(x) \approx \sum_{i=1}^m a_i \phi_i(x) \quad \text{if } m \rightarrow \infty, \\ \|f - \sum a_i \phi_i\| \rightarrow 0$$

$$a_i = \int f(x) \phi_i(x) dx$$

Ex: $\{\phi_i\}$: Legendre polynomial
Fourier basis

$f(x_1, x_2)$, 2-dim function $\{\phi_{i_1}(x_1)\}$, $\{\phi_{i_2}(x_2)\}$.

$$f(x_1, x_2) \approx \sum_{i_1, i_2=1}^m A_{i_1, i_2} \phi_{i_1}(x_1) \phi_{i_2}(x_2).$$

$A \in \mathbb{R}^{m \times m}$ coefficient matrix (sketching
↓
low-rank approx)

$$\underline{A_{i_1, i_2}} = \iint f(x_1, x_2) \phi_{i_1}(x_1) \phi_{i_2}(x_2) dx_1 dx_2$$

$f(x_1, \dots, x_d)$ d-dim function $\{\phi_{i_1}(x_1)\}$
 $\{\phi_{i_d}(x_d)\}$.

$$f(x_1, \dots, x_d) \approx \sum_{i_1, \dots, i_d=1}^m A_{i_1, i_2, \dots, i_d} \phi_{i_1}(x_1) \dots \phi_{i_d}(x_d)$$

$$A \in \mathbb{R}^{m \times m \times \dots \times m} \quad \underline{m}^d \text{ } \frac{d\text{-dim}}{\text{tensor}}$$

$$A_{i_1, \dots, i_d} = \int \dots \int f(x_1, \dots, x_d) \phi_{i_1}(x_1) \dots \phi_{i_d}(x_d) dx_1 \dots dx_d$$

2.2 Density Estimation. $f \xrightarrow{\text{prob}}$ density function p (pdf)

Given $\{X_i\}_{i=1}^N$ i.i.d. p^* , N = sample-size

Empirical density $\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X_i}(x)$
 $\hat{p} \rightarrow$ a random function based on X_i

$$\int \delta_{X_i}(x) dx = 1$$

Task: use some methods to get approximation of p^* based on \hat{p}

method
2.3 kernel density estimation (KDE):

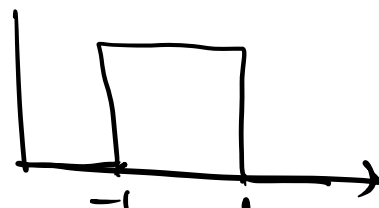
$$\hat{P}_h(x) = \frac{1}{N} \sum_{i=1}^N K_h(x - X_i) \leftarrow \text{1d.}$$

$K_h(\cdot) = \frac{1}{h} K\left(\frac{\cdot}{h}\right)$, $K(\cdot)$ - kernel function.

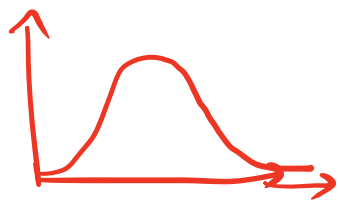
h - bandwidth.

Ex: step function $k(x) = \frac{1}{2} I(x)$

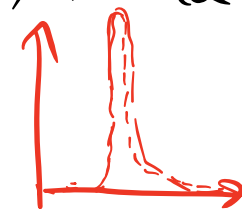
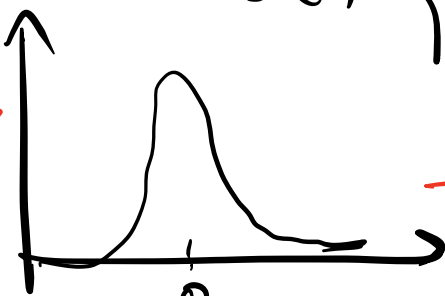
Gaussian kernel: $k(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$



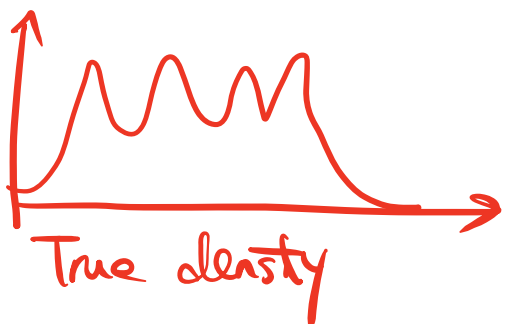
$$I(x) = \begin{cases} 1, & x \in [-1, 1] \\ 0, & x \text{ else} \end{cases}$$



h large



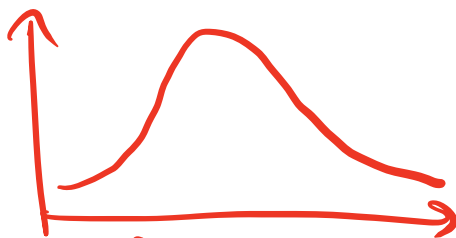
h small



True density



small- h .



Large- h

high-d

$$\hat{P}_h(x) = \frac{1}{N} \sum_{i=1}^N \frac{1}{h^d} k\left(\frac{\|x - X_i\|}{h}\right)$$

Error: bias error / variance error.

$$\mathbb{E} \left\| \hat{P}_h - p^* \right\|$$

$$\mathbb{E} \left\| \hat{P}_h - \mathbb{E} \hat{P}_h \right\|^2$$

suppose p^* belongs to Holder class.

$$\Sigma(\beta, L) = \left\{ g: |D^s g(x) - D^s g(y)| \leq L \|x - y\|, \right. \\ \left. \text{for all } s \text{ s.t. } |s| \leq \beta, \forall x, y \right\}$$

$$D^s = \frac{\partial^{s_1 + \dots + s_d}}{\partial x_1^{s_1} \dots \partial x_d^{s_d}}$$

$$\mathbb{E} \|\hat{p}_h - p^*\|^2 \leq \text{bias}^2 + \text{Variance} \leq ch^\beta + \frac{C}{Nh^d}$$

$$h \approx N^{-\frac{1}{2\beta+d}}$$

$$\approx \left(\frac{1}{N} \right)^{\frac{2\beta}{2\beta+d}}$$

$$d \nearrow, \frac{2\beta}{2\beta+d} \rightarrow 0$$

curse of dimensionality

Tensor. $\underline{n}^d \leftarrow$ low-rank.

keep error $\leq \varepsilon$ tolerance

$$\left(\frac{1}{N} \right)^{\frac{2\beta}{2\beta+d}} \leq \varepsilon$$

$$\frac{1}{\varepsilon} \leq N^{\frac{2\beta}{2\beta+d}}$$

$$N \geq \left(\frac{1}{\varepsilon} \right)^{\left(1 + \frac{d}{2\beta}\right)}$$

$d \nearrow, N \rightarrow \text{exp in } d, \text{ for fixed } \varepsilon.$

3. **Assumption 1:** $x \in \Omega_1 \subset \mathbb{R}^{d_1}$, $y \in \Omega_2 \subset \mathbb{R}^{d_2}$

let $A^*(x, y) = \Omega_1 \times \Omega_2 \rightarrow \mathbb{R}$ be a generic population function with $\|A^*\|_{L_2(\Omega_1 \times \Omega_2)} < \infty$.

$$\text{Assume } A^*(x, y) = \sum_{\rho=1}^r \sigma_{\rho} \Phi_{\rho}^*(x) \Psi_{\rho}^*(y)$$

r is const, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$

$\{\Phi_{\rho}^*(x)\}_{\rho=1}^r$, $\{\Psi_{\rho}^*(y)\}_{\rho=1}^r$ orthonormal function.

Ex: Mean-field model

$$p^*(x, y) = p_1^*(x) \cdot p_2^*(y)$$

$$p^*(z_1, \dots, z_d) = p_1^*(z_1) p_2^*(z_2) \cdot \dots \cdot p_d^*(z_d)$$

Extension: mixed mean-field.

$$p^*(z_1, \dots, z_d) = \sum_{\rho=1}^r \tau_{\rho} p_{\rho,1}^*(z_1) p_{\rho,2}^*(z_2) \cdot \dots \cdot p_{\rho,d}^*(z_d)$$

$$\text{span} \{\Phi_{\rho}^*(x)\} = \text{Span} \{p_{\rho,1}^*(z_1)\}, \text{Span} \{\Psi_{\rho}^*(y)\} \\ = \text{Span} \{p_{\rho,2}^*(z_2) \cdot \dots \cdot p_{\rho,d}^*(z_d)\}$$

rank at most r

How to get P? projection. Ex: Legendre poly basis.

sketching: matrix $S \rightarrow$ function $\{w_\eta(y)\}$
 space $L = \text{Span}\{w_\eta(y)\}_{\eta=1}^{\dim(L)}$

① sketching stage: $\int_{\Omega_2} \hat{A}(x,y) w_\eta(y) dy$ (integral)
 a function solely depends on x
 $[BS(:,1), BS(:,2), \dots, BS(:,k)]$

② Estimation Stage: space $M = \text{Span}\{v_\mu(x)\}_{\mu=1}^{\dim(M)}$

$\tilde{f}_\eta(x) = \arg \min_{f \in M} \left\| \int_{\Omega_2} \hat{A}(x,y) w_\eta(y) dy - f(x) \right\|_{L_2(\Omega_1)}^2$
 basis function.

$\tilde{f}_\eta(x) = \sum_{\mu=1}^{\dim(M)} B_{\mu,\eta} v_\mu(x)$

coefficient matrix $B_{\mu,\eta} = \left\langle \int_{\Omega_2} \hat{A}(x,y) w_\eta(y) dy, v_\mu(x) \right\rangle$
 $B \in \mathbb{R}^{\dim(M) \times \dim(L)}$
 $= \iint \hat{A}(x,y) v_\mu(x) w_\eta(y) dx dy$

③ Orthogonalization do QR/SVD $\rightarrow \underline{\underline{B}}$

$$[U, \Sigma, V] = \text{svd}(B)$$

$\mathbb{R}^{\text{dim}(M) \times r}$ (truncate at rank- r).

$$\hat{\Phi}_p(x) = \sum_{\mu=1}^{\text{dim}(M)} v_{\mu}(x) U_{\mu,p}, \quad p=1, \dots, r$$

projection operator $P_x(x, x') = \sum_{p=1}^r \hat{\Phi}_p(x) \hat{\Phi}_p(x')$

$$\begin{array}{c} \downarrow \\ U U^T \\ \Downarrow \\ U(x) U^T(x') \end{array}$$

Algorithm:

compute coefficient B :

$$B = \frac{1}{N} \sum_{i=1}^N v_{\mu}(x_i) w_{\nu}(y_i).$$

$$\text{when } \hat{A}(x, y) = \frac{1}{N} \sum_{i=1}^N \delta(x_i, y_i)(x, y)$$

Function estimation. (2-d function).

$$\dim(L_2) \rightarrow \frac{1}{G_r}$$

$$\dim(M_2) \rightarrow$$

$$\textcircled{1} \left\{ \hat{\Phi}_p(x) \right\}_{p=1}^r$$

sketching space $\left\{ W_\eta(x) \right\}_{\eta=1}^{\dim(L_2)}$
 estimation space $\left\{ V_\mu(y) \right\}_{\mu=1}^{\dim(M_2)}$

$$\hat{A}(x, y) \rightarrow$$

②

$$\text{Algorithm} \rightarrow \left\{ \hat{\Phi}_p(y) \right\}_{p=1}^r$$

\hat{A} is matrix
 ① projection: row $U_1 U_1^T \hat{A}$
 ② projection: column $\hat{A} U_2 U_2^T$
 ③ Together $\tilde{A} = U_1 U_1^T \hat{A} U_2 U_2^T$

③ efficient core.

$$G_{p_1, p_2} = \iint \hat{A}(x, y) \hat{\Phi}_{p_1}(x) \hat{\Phi}_{p_2}(y) dx dy$$

$\xrightarrow{\text{density}} \frac{1}{N} \sum_{i=1}^N \hat{\Phi}_{p_1}(x_i) \hat{\Phi}_{p_2}(y_i)$

Represent function

$$\tilde{A}(x, y) = \sum_{p_1, p_2=1}^r G_{p_1, p_2} \hat{\Phi}_{p_1}(x) \hat{\Phi}_{p_2}(y)$$

$$\tilde{A} = \hat{A} \underset{=}{X_1} \underset{=}{P_1} X_2 \underset{=}{P_2}$$

Statistical analysis · $x \in \mathbb{R}^{d_1}, y \in \mathbb{R}^{d_2}$

Assumption 2: let M, L with $\dim(M) = m^{d_1}$,
 $\dim(L) = l^{d_2}$

$$\|A^* - A^* x x^T P_M x y^T P_L\|_{L_2}^2 = O(m^{-2\alpha} + l^{-2\alpha})$$

α : parameter of space (containing A^*).

Ex: A^* in Sobolev space W_2^α : up to α derivative is bounded in L_2^{non}

$$\|A^*\|_{W_2^\alpha}^2 = \left(\sum_{|S|=0}^{\alpha} \|D^S A^*\|_2^2 \right) < \infty$$

$D^S = \frac{\partial^{s_1 + \dots + s_d}}{\partial x_1^{s_1} \dots \partial x_d^{s_d}}$

Intuition: basis in M : polynomial · $(s_1 + \dots + s_d = |s|) x \in \mathbb{R}^d$

$\dim(M)$: how many orthogonal polynomials

→ degree of polynomial

to each dimensional order of polynomial

$$\dim(M) = m^{d_1}, x \in \mathbb{R}^{d_1}$$

Error ↓

Ex: $\alpha = 2$

$$W_2^2, f \in W_2^2, \|f\|_{L_2} < \infty, \|f'\|_{L_2} < \infty, \|f''\|_{L_2} < \infty \text{ for } \forall x$$

Assumption 3: for any test function $u(x, y)$,

$$\textcircled{1} \mathbb{E}_{\text{samples}} \langle \hat{A}, u \rangle = \langle A^*, u \rangle \quad \textcircled{2} \sup_{\|u\|_{L_2} \leq 1} \text{Var}[\langle \hat{A}, u \rangle] = O\left(\frac{1}{N}\right)$$

1st moment
sample-size
2nd moment

$$\left\{ \begin{aligned} \langle \hat{A}, u \rangle &= \iint \hat{A}(x, y) u(x, y) dx dy \\ \langle A^*, u \rangle &= \iint A^*(x, y) u(x, y) dx dy \end{aligned} \right.$$

$$\hat{A} = \frac{1}{N} \sum \delta_i, \quad \textcircled{2} = \sup |A^*(x, y)| < \infty$$

Theorem 1: (informal) suppose above assumptions hold,

$$\|\tilde{A} - A^*\|_{L_2}^2 = O\left(\frac{1}{N^{\frac{1}{2\alpha + d_1}}} + \frac{1}{N^{\frac{1}{2\alpha + d_2}}}\right)$$

$x \in \mathbb{R}^{d_1}$
 $y \in \mathbb{R}^{d_2}$

with choice $m_1 \approx N^{\frac{1}{2\alpha + d_1}}, m_2 \approx N^{\frac{1}{2\alpha + d_2}}$

sketching size \downarrow $\dim(M_1) \rightarrow$ variable x variable y .

$$u \approx l_2 \approx C \epsilon_r^{-\frac{1}{\alpha}}$$

\rightarrow r -th singular value of A^*

$$\text{KDE: } O\left(N^{\frac{1}{2\alpha + d_1 + d_2}}\right)$$

bias-variance trade-off.

Remark: slightly improvement for 2-d function.
 huge improvement for high-d function

3.3 high-d function estimation.

Assumption 4: (modification of low-rank..)

$$A^*(x_1, \dots, x_d) = \sum_{\rho=1}^{r_j} \sigma_{j,\rho} \Phi_{j,\rho}^*(x_j) \Psi_{j,\rho}^*(x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_d)$$

for $j=1, \dots, d$

sketching space: num of basis for x_j

$$M_j = \text{span} \{ \phi_{\mu}(x_j) \}_{\mu=1}^m \quad (d-1)$$

$$L_j = \text{span} \{ \phi_{\eta_1}(x_1) \dots \phi_{\eta_{j-1}}(x_{j-1}) \phi_{\eta_{j+1}}(x_{j+1}) \dots \phi_{\eta_d}(x_d) \}_{\eta_1, \dots, \eta_d=1}^{l_j}$$

$\dim(L_j) = l_j^{d-1}$

$\left. \begin{array}{l} m \text{ could be large} \\ l_j \text{ small} \end{array} \right\}$

$l_j \rightarrow$ sketching.

Thm 2: (informal) Suppose above assumption holds,
choose $m \approx N^{\frac{1}{2\alpha+1}}$, $G_j = C G_j r_j^{-\frac{1}{\alpha}}$,

then

$$\|\tilde{A} - A^*\|_{L_2}^2 = O\left(\frac{1}{N^{\frac{1}{2\alpha(2\alpha+1)}}}\right)$$

Remark 1: intuitive: $M_j \rightarrow \mathcal{X}_j \in \mathbb{R}^d$.
 $d-1$ dimensions are sketched. $N^{\frac{1}{2\alpha(2\alpha+1)}}$

Remark 2: we sketch $d-1$ dimensions

$$\dim(L_j) = L^{d-1}$$

image L is small when A^* is good.