

Linear Bandits

Least squares review

Suppose I measure $x_1, \dots, x_T \in \mathbb{R}^d$ then observe

that $y_t = \langle x_t, \theta_* \rangle + \zeta_t$ where ζ_t is mean-zero σ^2 sub-Gaussian

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{t=1}^T (y_t - \langle x_t, \theta \rangle)^2$$

$$= \underset{\theta}{\operatorname{argmin}} \|y - X\theta\|_2^2$$

$$= (X^T X)^{-1} X^T y$$

$$= (X^T X)^{-1} X^T (X\theta_* + \zeta)$$

$$= \theta_* + (X^T X)^{-1} X^T \zeta$$

$$X = \begin{bmatrix} -x_1^T \\ \vdots \\ -x_T^T \end{bmatrix} \in \mathbb{R}^{T \times d}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_T \end{bmatrix} \in \mathbb{R}^T$$

$$W = X(X^T X)^{-1} \quad \text{then} \quad \hat{\theta} = \theta_* + W^T \zeta$$

$$\mathbb{E}[(\hat{\theta} - \theta_*)(\hat{\theta} - \theta_*)^T] = \mathbb{E}[(X^T X)^{-1} X^T \zeta \zeta^T X (X^T X)^{-1}]$$

$$= (X^T X)^{-1} X^T \mathbb{E}[\zeta \zeta^T] X (X^T X)^{-1}$$

$$\leq \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1}$$

Let $z \in \mathbb{R}^d$ be some new point.

$$\mathbb{E}[\langle z, \hat{\theta} \rangle] = \langle z, \theta_* \rangle$$

$$\mathbb{V}(\langle z, \hat{\theta} \rangle) = \mathbb{E}[\langle z, \hat{\theta} - \theta_* \rangle^2]$$

$$= \mathbb{E}[z^T (\hat{\theta} - \theta_*)(\hat{\theta} - \theta_*)^T z]$$

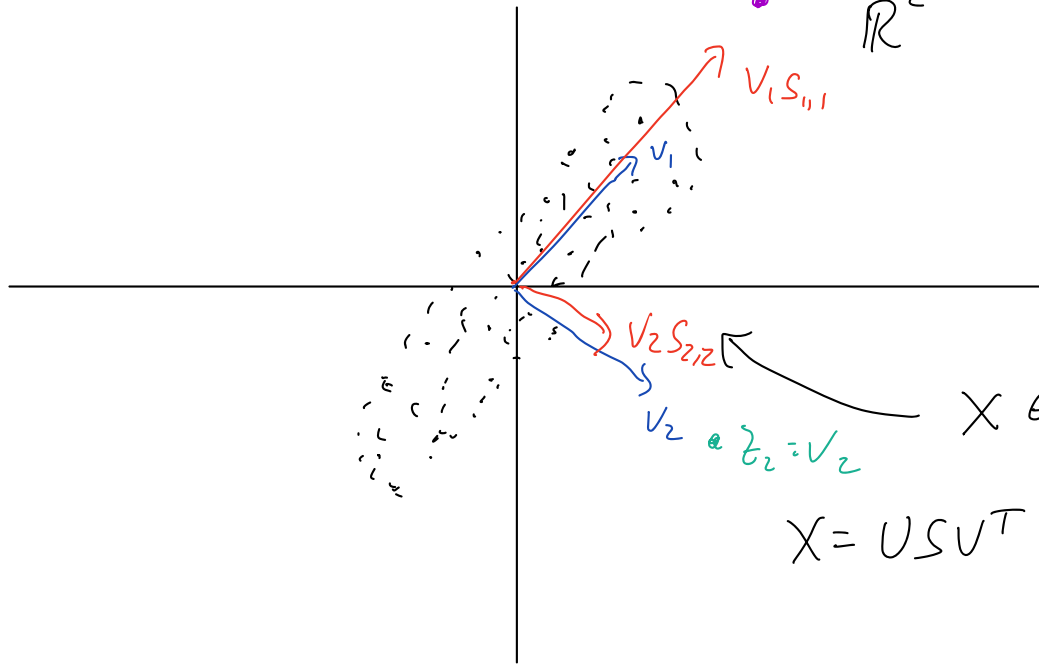
For any PSD $A \in \mathbb{R}^{d \times d}$

Define $\|x\|_A^2 = x^T A x$.

$$\leq \sigma^2 z^T (X^T X)^{-1} z$$

$$= \sigma^2 \|z\|_{(X^T X)^{-1}}^2$$

$$z_1 = v_1 \quad \mathbb{R}^2$$



$$X \in \mathbb{R}^{T \times 2}$$

$$X = U S V^T$$

$$U \in \mathbb{R}^{T \times 2}$$

$$S \in \mathbb{R}^{2 \times 2}$$

$$V \in \mathbb{R}^{2 \times 2}$$

$$U = \begin{bmatrix} -u_1^T \\ \vdots \\ -u_1^T \end{bmatrix}$$

$$U^T U = I_2$$

$$V^T V = I_2$$

$$V = \begin{bmatrix} | & | \\ v_1 & v_2 \\ | & | \end{bmatrix}$$

$$x_i = u_i^T S V$$

$$= u_{i,1} s_{11} v_1 + u_{i,2} s_{22} v_2$$

$$(X^T X)^{-1} = (V S U^T U S V^T)^{-1}$$

$$= (V S^2 V^T)^{-1}$$

$$= V S^{-2} V^T$$

V is orthonormal

$$V^{-1} = V^T$$

If $z_1 = v_1$ then

$$\begin{aligned} z_1^T (X^T X)^{-1} z_1 &= v_1^T V S^{-2} V^T v_1 \\ &= e_1^T \bar{S}^{-2} e_1 \\ &= \frac{1}{s_{1,1}^2} \end{aligned}$$

If $z_2 = v_2$

$$z_2^T (X^T X)^{-1} z_2 = \frac{1}{s_{2,2}^2}$$

How do I choose $x_1, \dots, x_T \in \mathcal{X}$ to minimize $\mathbb{E}[\langle z, \hat{\theta} - \theta_* \rangle^2]$?

Suppose I have a pool of potential measurements $\mathcal{X} \subset \mathbb{R}^d$.

For simplicity assume $|\mathcal{X}| = n$.

For any $X^T X = \sum_{i=1}^T x_i x_i^T \exists \lambda \in \Delta_{\mathcal{X}} = \{\omega_i \geq 0 : \sum_{i=1}^n \omega_i = 1\}$ s.t. $X^T X = T \sum_{x \in \mathcal{X}} \lambda_x x x^T$.

$$A(\lambda) = \sum_{x \in \mathcal{X}} \lambda_x x x^T$$

(Always exists by $\lambda_x := \frac{\sum_{i=1}^T \mathbb{1}\{x_i = x\}}{T}$)
=: $T \cdot A(\lambda)$

$$\begin{aligned} \mathbb{E}[\langle z, \hat{\theta} - \theta_* \rangle^2] &\leq \sigma^2 z^T (X^T X)^{-1} z \\ &= \sigma^2 \frac{z^T A(\lambda)^{-1} z}{T} \end{aligned}$$

$$A(\lambda) = \mathbb{E}_{x \sim \lambda} [x x^T]$$

Idea: Choose $\lambda \in \Delta_{\mathcal{X}}$ to minimize RHS! Before seeing any data then choose x_1, \dots, x_T "in proportion" to $\lambda \in \Delta_{\mathcal{X}}$.

A-optimality

$$\begin{aligned}\sum_{i=1}^d \mathbb{E}[\langle e_i, \hat{\theta} - \theta_* \rangle^2] &= \sum_{i=1}^d \mathbb{E}[(\hat{\theta} - \theta_*)^T e_i e_i^T (\hat{\theta} - \theta_*)] \\ &= \mathbb{E}[(\hat{\theta} - \theta_*)^T (\hat{\theta} - \theta_*)] \\ &= \mathbb{E}[\|\hat{\theta} - \theta_*\|_2^2]\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\|\hat{\theta} - \theta_*\|_2^2] &= \sum_{i=1}^d \mathbb{E}[\langle e_i, \hat{\theta} - \theta_* \rangle^2] & [c_{ij}] &= \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases} \\ &\leq \sum_{i=1}^d \frac{\sigma^2 e_i^T A(\lambda)^{-1} e_i}{T} \\ &= \frac{\sigma^2}{T} \cdot \text{Trace}(A(\lambda)^{-1}).\end{aligned}$$

A-optimal design: $\underset{\lambda \in \Delta_K}{\text{argmin}} \text{Trace}(A(\lambda)^{-1})$.

E-optimality

$$\begin{aligned}\max_{u: \|u\|_2 \leq 1} \mathbb{E}[\langle u, \hat{\theta} - \theta_* \rangle^2] &\leq \frac{\sigma^2}{T} \max_{u: \|u\|_2 \leq 1} u^T A(\lambda)^{-1} u \\ &= \frac{\sigma^2}{T} \lambda_{\max}(A(\lambda)^{-1}) \\ &= \frac{\sigma^2}{T} \cdot \frac{1}{\lambda_{\min}(A(\lambda))}\end{aligned}$$

(e_i, v_i) are eigen pairs

$$A v_i = e_i v_i$$

$$f_E(\lambda) = \lambda_{\max}(A(\lambda)^{-1})$$

D-optimality

$$g_D(\lambda) = \log |A(\lambda)|$$

λ^* is D-optimal

if $\lambda^* = \underset{\lambda}{\text{argmax}} g_D(\lambda)$.

D-optimal maximizes product of eigs of $A(\lambda)$.

$\mathcal{N}(\mu, \Sigma)$ has density $\frac{1}{(2\pi|\Sigma|)^{d/2}} \exp(-(x-\mu)^T \Sigma^{-1} (x-\mu))$.

G-optimal

$$f_G(\lambda) = \max_{x \in \mathcal{X}} \|x\|_{A(\lambda)^{-1}}^2 = \max_{x \in \mathcal{X}} x^T A(\lambda)^{-1} x$$

G-opt minimizes $\max_{x \in \mathcal{X}} \mathbb{E}[(x, \hat{\theta} - \theta_*)^2]$.

Kiefer-Wolfowitz Equivalence Theorem

Let $\lambda_d = \operatorname{argmax}_{\lambda} g_D(\lambda)$. Then

$$g_D(\lambda_d) = \max_{\lambda} g_D(\lambda)$$

$$f_G(\lambda_d) = \min_{\lambda} f_G(\lambda) = d.$$

Intuition for $f_G(\lambda_d) = d$:

$$\begin{aligned} \max_x \|x\|_{A(\lambda)^{-1}}^2 &\geq \sum_x \lambda_x \|x\|_{A(\lambda)^{-1}}^2 \\ &= \sum_x \lambda_x x^T A(\lambda)^{-1} x \\ &= \sum_x \lambda_x \operatorname{Tr}(x^T A(\lambda)^{-1} x) \\ &= \sum_x \lambda_x \operatorname{Tr}(A(\lambda)^{-1} x x^T) \\ &= \operatorname{Tr}(A(\lambda)^{-1} A(\lambda)) \\ &= d \end{aligned}$$

How do I turn $\lambda \in \Delta_x$ into x_1, \dots, x_T ?

If $T\lambda \in \mathbb{N}^n$ then take $Tx = \lambda_x T$.

1) One idea is to $x_i \stackrel{iid}{\sim} \lambda$, $i=1, \dots, T$.

$$\Rightarrow \sum_{i=1}^T x_i x_i^T = A(\lambda) \cdot T$$

Problem:

$$A(\lambda)^{-1} = \left(\mathbb{E} \left[\sum_{i=1}^T x_i x_i^T \right] \right)^{-1} \neq \mathbb{E} \left[\left(\sum_{i=1}^T x_i x_i^T \right)^{-1} \right]$$

Not going to work

2) Another idea is to take $Tx = \lceil \lambda_x T \rceil$

Problem: might go over budget: $\sum_x Tx \gg T$.

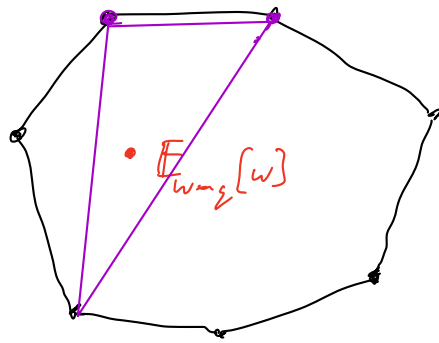
Fact For every $\lambda \in \Delta_x$, $\exists \lambda' \in \Delta_x$ s.t. $\|\lambda'\|_0 \leq \frac{d(d+1)}{2}$
that satisfies $A(\lambda) = A(\lambda')$.

Why? $A(\lambda) = \sum_x \lambda_x \cdot xx^T$.

Caratheodory's Theorem Let $\Omega \subset \mathbb{R}^p$ and let $g \in \Delta_\Omega$
Then $\exists g'$: $\mathbb{E}_{w \sim g} [w] = \mathbb{E}_{w \sim g'} [w]$, and $\|g'\|_0 \leq p+1$.

Proof by picture

Ω



G-optimal

Procedure | Input: \mathcal{X} , budget T

1.) Solve $\tilde{\lambda}_* = \underset{\lambda \in \Delta_{\mathcal{X}}}{\operatorname{argmin}} \max_{x \in \mathcal{X}} \|x\|_{A(\lambda)}^2$

2.) Find a $\frac{d(d+1)}{2}$ sparse solution λ_* : $A(\lambda_*) = A(\tilde{\lambda}_*)$.

3.) Set $T_x = \lceil \lambda_{*,x} T \rceil$ times, observe $y(x) := \langle x, \theta_* \rangle + z$
 T_x times to obtain

4.) Compute $\hat{\theta} = (X^T X)^{-1} X^T y$ $\{z_i\}_{i=1}^{\sum_x T_x}$

Lemma $\sum_x T_x = \sum_x \lceil \lambda_{*,x} T \rceil \leq \frac{d(d+1)}{2} + \sum_x \lambda_{*,x} T = \frac{d(d+1)}{2} + T$

With probability at most $1-\delta$ we have for all $x \in \mathcal{X}$

$$|\langle x, \hat{\theta} - \theta_* \rangle| \leq \|x\|_{A(\lambda_*)} \sqrt{\frac{2d\sigma^2 \log(2|x|/\delta)}{T}} = \sqrt{\frac{2d\sigma^2 \log(2|x|/\delta)}{T}}$$

For any z we have $z^T (\hat{\theta} - \theta_*) = z^T (X^T X)^{-1} X^T z =: w^T z$

$$\mathbb{E}[w^T z] = 0, \quad \mathbb{E}[(w^T z)^2] = w^T \mathbb{E}[z z^T] w = \leq \sigma^2 \|w\|_2^2$$

$$\begin{aligned} \mathbb{E}[\exp(\lambda w^T z)] &= \mathbb{E}[\exp(\lambda \sum_t w_t z_t)] && z_t \text{ is } \mathcal{N}(0, \sigma^2) \text{-Gaussian} \\ &= \mathbb{E}[\prod_t \exp(\lambda w_t z_t)] \end{aligned}$$

$$= \prod_t \mathbb{E}[\exp(\lambda w_t z_t)]$$

$$\leq \prod_t \exp(\lambda^2 w_t^2 \sigma^2 / 2)$$

$$= \exp\left(\frac{\lambda^2 \sigma^2}{2} \sum_t w_t^2\right)$$

$$= \exp\left(\frac{\lambda \sigma^2}{2} \|w\|_2^2\right)$$

For any $\lambda' \in \mathbb{R}$ we have $\mathbb{E}[\exp(\lambda' z_t)] \leq \exp(\lambda'^2 \sigma^2 / 2)$

$$\|w\|_2^2 = w^T w = z^T (X^T X)^{-1} X^T X (X^T X)^{-1} z = z^T (X^T X)^{-1} z$$

$\Rightarrow (\hat{\theta} - \theta_*)^T z$ is $(\sigma^2 \|z\|_{(X^T X)^{-1}}^2)$ -sub-Gaussian.

For any $z \in \mathcal{X}$

$$\mathbb{P}((\hat{\theta} - \theta_*)^T z > \varepsilon) \leq \exp\left(-\frac{\varepsilon^2}{2\sigma^2 \|z\|_{(X^T X)^{-1}}^2}\right) \quad (\text{Chernoff})$$

$$\leq \exp\left(-\frac{T\varepsilon^2}{2d\sigma^2}\right)$$

$$\|z\|_{(X^T X)^{-1}}^2 = z^T (X^T X)^{-1} z = z^T \left(\sum_x \lambda_x^{-1} x x^T\right)^{-1} z$$

$$\leq z^T \left(\sum_x \lambda_x^{-1} x x^T\right)^{-1} z$$

$$= \frac{z^T \Lambda(\lambda)^{-1} z}{T} = \frac{\|z\|_{\Lambda(\lambda)^{-1}}^2}{T} = \frac{d}{T}$$

$$\mathbb{P}\left(\bigcup_{z \in \mathcal{X}} \{z^T (\hat{\theta} - \theta_*) > \varepsilon\}\right) \leq 2 \sum_{z \in \mathcal{X}} \mathbb{P}(z^T (\hat{\theta} - \theta_*) > \varepsilon)$$

$$\leq 2|\mathcal{X}| \exp\left(-\frac{T\varepsilon^2}{2d\sigma^2}\right). \quad \text{Set } \delta \text{ solve for } \varepsilon. \quad \square$$