

Fix  $x_1, \dots, x_n \in \mathbb{R}^d$  and  $\theta_* \in \mathbb{R}^d$ .

known  
↓

unknown  
↓

At  $i$  Observe

$$y_i = x_i^T \theta_* + \varepsilon_i, \quad \{\varepsilon_i\} \text{ are IID mean-zero 1-sub-Gaussian R.V.}$$

Least squares estimate:

$$\begin{aligned} \hat{\theta} &= \arg \min_{\theta} \sum_{i=1}^n (\langle x_i, \theta \rangle - y_i)^2 \\ &= \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i \\ &= \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} \left( \sum_{i=1}^n x_i x_i^T \right) \theta_* + \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i \varepsilon_i \\ &= \theta_* + \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i \varepsilon_i \end{aligned}$$

Observe

$$\{(x_i, y_i)\}_{i=1}^n$$

$$\mathbb{E}[\hat{\theta}] = \theta_*$$

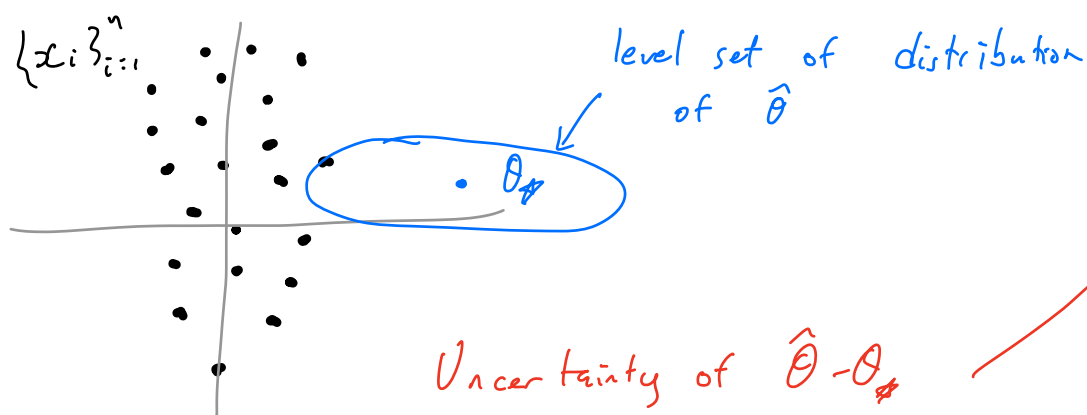
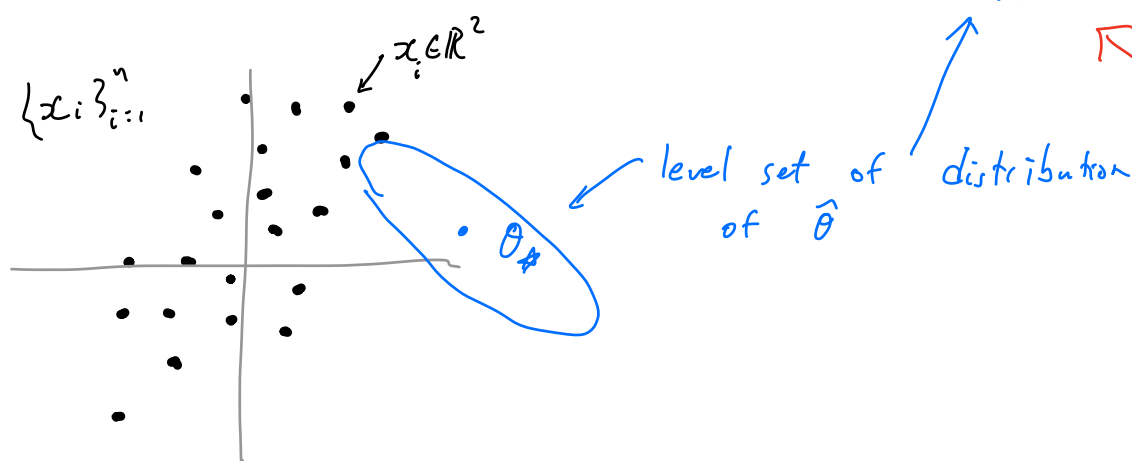
$$\text{Cov}(\hat{\theta}) = \left( \sum_{i=1}^n x_i x_i^T \right)^{-1}$$

Fix any  $z \in \mathbb{R}^d$ . With probability at least  $1-\delta$

$$\langle z, \hat{\theta} - \theta_* \rangle \leq \|z\|_{\left( \sum_{i=1}^n x_i x_i^T \right)^{-1}} \sqrt{2 \log(1/\delta)}. \quad (\|x\|_A^2 = x^T A x)$$

Proof: Show  $\langle z, \hat{\theta} - \theta_* \rangle$  is  $\|z\|_{\left( \sum_{i=1}^n x_i x_i^T \right)^{-1}}$ -sub-Gaussian. Apply Chernoff bound.

Intuition: If  $\varepsilon_i \sim \mathcal{N}(0, 1)$  then  $\tilde{\theta} \sim \mathcal{N}(\theta_*, (\sum_{i=1}^n x_i x_i^T)^{-1})$



Uncertainty of  $\hat{\theta} - \theta_*$   
depends only on  $\{x_i\}_{i=1}^n$  and not  $\theta_*$  or  $\{q_i\}$

Experimental Design: Exploits this observation  
to choose  $\{x_i\}$  in order to obtain a  
desired covariance shape.

Given choice of  $\{x_i\}_{i=1}^n$ , and then observe  $y_i = \langle x_i, \theta_* \rangle + \varepsilon_i$   
results in  $\tilde{\theta}$  w/  $\mathbb{E}[\tilde{\theta}] = \theta_*$   $\text{Cov}(\tilde{\theta}) = (\sum_i x_i x_i^T)^{-1}$

Given some pool of points  $\mathcal{X} \subset \mathbb{R}^d$

choose  $\{x_i\}_{i=1}^n$

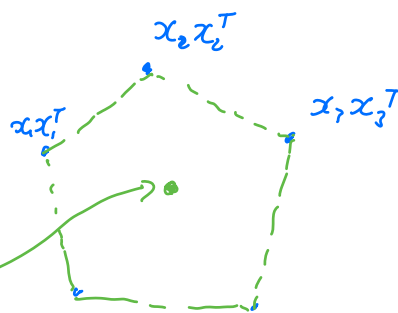
For every choice of points  $\{x_i\}_{i=1}^n \subset \mathcal{X}$ ,  $\exists \lambda \in \Delta_{\mathcal{X}}$

$$\Delta_{\mathcal{X}} = \left\{ p \in \mathbb{R}^{|\mathcal{X}|} : \sum_{x \in \mathcal{X}} p_x = 1, p_x \geq 0 \ \forall x \in \mathcal{X} \right\}$$

s.t. 
$$\sum_{i=1}^n x_i x_i^T = n \sum_{x \in \mathcal{X}} \lambda_x x x^T$$

why? 
$$\lambda_x = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = x\},$$

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



Exp. Design Objectives

- A-optimality  $f_A(\lambda) = \text{Tr}(A(\lambda)^{-1})$

$$\mathbb{E}[\|\hat{\theta} - \theta_*\|_2^2] = \text{Tr}\left(\left(\sum_{i=1}^n x_i x_i^T\right)^{-1}\right)$$

//

minimize average error over all directions"

- E-optimality  $f_E(\lambda) = \sup_{u: \|u\|_2 \leq 1} u^T A(\lambda)^{-1} u$

$$\sup_{u: \|u\|_2 \leq 1} \mathbb{E} \left[ \langle u, \hat{\theta} - \theta_* \rangle^2 \right] = \sup_u u^T \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} u$$

"minimize worst-case error direction"

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

$$\begin{aligned} \max_{x \in \mathcal{X}} \mathbb{E} \left[ \langle x, \hat{\theta} - \theta_* \rangle^2 \right] &= \max_{x \in \mathcal{X}} x^T \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} x \\ &= \max_{x \in \mathcal{X}} \|x\|_{\left( \sum_{i=1}^n x_i x_i^T \right)^{-1}}^2 \end{aligned}$$

- D-optimality  $f_D(\lambda) = -\log(|A(\lambda)|)$   
 $= \log(|A(\lambda)^{-1}|)$

Recall: if  $Z \sim \mathcal{N}(0, \Sigma)$

then entropy of  $Z = \frac{1}{2} \log(2\pi e |\Sigma|)$

"minimize entropy of  $\hat{\theta}$ "

Lemma (Kiefer-Wolfowitz 1960)

For any  $X \subset \mathbb{R}^d : \dim(\text{span}(X)) = d$ , there

exists a  $\lambda^* \in \Delta_X$  such that

$$- \min_{\lambda} f_D(\lambda) = f_D(\lambda^*)$$

$$- \min_{\lambda} f_G(\lambda) = f_G(\lambda^*)$$

$$- f_G(\lambda^*) = d$$

$$- \text{support}(\lambda^*) \leq \frac{(d+1)d}{2}.$$

Caratheodory  
Theorem.

$$\min_{\lambda} f_G(\lambda) = \min_{\lambda} \max_{x \in X} x^T A(\lambda)^{-1} x$$

$$\geq \min_{\lambda} \sum_{x \in X} \lambda_x x^T A(\lambda)^{-1} x$$

$$= \min_{\lambda} \text{Tr} \left( \underbrace{\sum_{x \in X} \lambda_x x x^T}_{= I} A(\lambda)^{-1} \right)$$

$x \in \mathbb{R}^d$  is

$s$ -sparse if

$|\{i \in [d] : x_i \neq 0\}| = s.$

$$= d$$

Proposition Fix  $X \subset \mathbb{R}^d$ . Assume  $Y = \langle X, \theta_* \rangle + \epsilon$  w/

Σ 1-sub-Gaussian. If  $\lambda^* = \min_{\lambda \in \Delta_K} f_G(\lambda)$  is a  $\frac{(d+1)d}{2}$ -sparse

solution, and we pull arm  $x \in \mathcal{X}$  exactly

$\lceil 3/\epsilon \rceil$  times, then w.p.  $\geq 1 - \delta$

$$|\langle x, \hat{\theta} - \theta_* \rangle| \leq \underbrace{\sqrt{\frac{2 d \log(2|x|/f)}{3}}}_{=\varepsilon} \quad \forall x \in \mathcal{X}$$

$$P\left(\bigcup_{x \in \mathcal{X}} \{ \langle x, \hat{\theta} - \theta_* \rangle > \varepsilon \} \right) \leq \exp\left(-\frac{\varepsilon^2}{2d}\right) \cdot 2|\mathcal{X}|$$

and  $\underbrace{\text{total \# pulls}}_{=n} \leq 3 + \frac{(d+1)d}{2}.$

From above, for any  $x \in \mathcal{X}$  w.p.  $\geq 1 - \delta'$

$$\langle x, \hat{\theta} - \theta_{\#} \rangle \leq \|x\|_{\left(\sum_{i=1}^n x_i x_i^T\right)^{-1}} \sqrt{2 \log(1/\delta')}$$

$$\sum_{i=1}^n x_i x_i^T$$

$$= \sum_{x \in \mathcal{X}} |\mathcal{S}_x^*| x x^T$$

$$\geq \sum_{x \in \mathcal{X}} \lambda_x x x^T$$

$$\leq \|x\| \left( 3 \sum_{x \in \mathcal{X}} \lambda_x^* x x^T \right)^{-1} \sqrt{2 \log(1/\delta')}$$

$$\leq \max_{x' \in \mathcal{X}} \|x'\| \left( \sum_{x \in \mathcal{X}} t_x^* x x^T \right)^{1/2} \sqrt{\frac{2 \log(1/\delta')}{3}}$$

 $A \succeq B$  if

$$A - B \succeq 0 \quad (\text{PSP})$$

$$= \sqrt{\frac{2d \log(1/\delta')}{3}}$$

$$\begin{aligned} \|x\|_{\left(\sum_{x' \in \mathcal{X}} \lambda_{x'} x' x'^T\right)^{-1}} &= \sqrt{x^T \left(\sum_{x' \in \mathcal{X}} \lambda_{x'} x' x'^T\right)^{-1} x} \\ &= \frac{1}{\sqrt{3}} \cdot \|x\|_{A(G^*)^{-1}} \end{aligned}$$

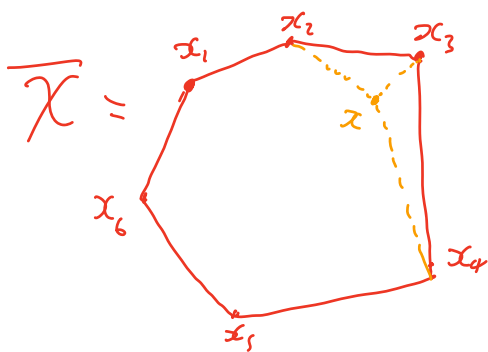
### Caratheodory

let  $x_1, \dots, x_n \in \mathbb{R}^d$  and define

$$\overline{\mathcal{X}} = \left\{ \sum_{i=1}^n p_i x_i : \sum_{i=1}^n p_i = 1, p_i \geq 0 \forall i \right\}$$

Then for any  $x \in \overline{\mathcal{X}}$ ,  $\exists p \in \Delta_n : \sum_{i=1}^n p_i x_i = x$

$$\|p\|_0 = d+1$$



**Input:** Finite set  $\mathcal{X} \subset \mathbb{R}^d$ , confidence level  $\delta \in (0, 1)$ .

Let  $\hat{\mathcal{X}}_1 \leftarrow \mathcal{X}, \ell \leftarrow 1$

**while**  $|\mathcal{X}_\ell| > 1$  **do**

Let  $\hat{\lambda}_\ell \in \Delta_{\mathcal{X}_\ell}$  be a  $\frac{d(d+1)}{2}$ -sparse minimizer of  $f(\lambda) = \max_{x \in \mathcal{X}_\ell} \|x\|^2_{(\sum_{x \in \mathcal{X}_\ell} \lambda_x x x^\top)^{-1}}$

$\epsilon_\ell = 2^{-\ell}, \tau_\ell = 2d\epsilon_\ell^{-2} \log(4\ell^2|\mathcal{X}|/\delta)$

Pull arm  $x \in \mathcal{X}_\ell$  exactly  $\lceil \hat{\lambda}_{\ell,x} \tau_\ell \rceil$  times and construct the least squares estimator  $\hat{\theta}_\ell$  using only the observations of this round

$\mathcal{X}_{\ell+1} \leftarrow \mathcal{X}_\ell \setminus \{x \in \mathcal{X}_\ell : \max_{x' \in \mathcal{X}_\ell} \langle x' - x, \hat{\theta}_\ell \rangle > 2\epsilon_\ell\}$

$\ell \leftarrow \ell + 1$

**Output:**  $\mathcal{X}_\ell$

$$\begin{aligned} \text{Regret} \quad R_T &= \max_{x \in \mathcal{X}} \mathbb{E} \left[ \sum_{t=1}^T \langle x, \theta_\# \rangle - \langle x_t, \theta_\# \rangle \right] \\ &\quad \uparrow \\ &\quad \text{arm algorithm played} \\ &\quad \text{at time } t \\ &= \sum_{x \in \mathcal{X}} \Delta_x \mathbb{E}[T_x] \end{aligned}$$

$$\Delta_x = \max_{x' \in \mathcal{X}} \langle x' - x, \theta_\# \rangle$$

$$T_x = \sum_{t=1}^T \mathbb{1}\{x = x_t\}$$

Assume  $\max_{x \in \mathcal{X}} |\langle x, \theta_\# \rangle| \leq 1$ .

Let  $x^\# = \arg\max_{x \in \mathcal{X}} \langle x, \theta_\# \rangle$

Lemma With prob at least  $1 - \delta$ , we have  $x^\# \in \mathcal{X}_\ell$

and  $\max_{x \in \mathcal{X}_\ell} \Delta_x \leq 8\epsilon_\ell$  for all  $\ell \in \mathbb{N}$ .

What does linear structure get us?

- If we ignored structure our elim. alg of last time has  $R_T \leq \sqrt{|\mathcal{X}| T \log(T)}$ .

- This algorithm that exploits linear structure,

we will show  $R_T \leq \sqrt{dT \log(|X|T)}$ .

I can "cover"  $X$  with just  $O(2^d)$  landmarks

$$\Rightarrow R_T \leq \sqrt{d^2 T \log(T)}$$

Proof. Consider  $V \subset X$  and event

$$\mathcal{C}_{x,l}^d(V) = \left\{ |\langle x, \hat{\theta}_l(V) - \theta^* \rangle| \leq \varepsilon_l \right\}$$

↑

The  $\hat{\theta}_l$  that would arise if  $X_l = V$ .

$$\mathbb{P} \left( \bigcup_{l=1}^{\infty} \bigcup_{x \in X_l} \mathcal{C}_{x,l}^d(X_l) \right) \leq \sum_{l=1}^{\infty} \mathbb{P} \left( \bigcup_{x \in X_l} \mathcal{C}_{x,l}^d(X_l) \right)$$

$$= \sum_{l=1}^{\infty} \sum_{V \subset X} \mathbb{P} \left( \left\{ \bigcup_{x \in V} \mathcal{C}_{x,l}^d(V) \right\} \cap \{X_l = V\} \right)$$

$$= \sum_{l=1}^{\infty} \sum_{V \subset X} \underbrace{\mathbb{P} \left( \bigcup_{x \in V} \mathcal{C}_{x,l}^d(V) \right)}_{\text{apply proposition}} \mathbb{P}(X_l = V)$$

$$\leq \sum_{l=1}^{\infty} \sum_{V \subset X} \underbrace{\frac{|V|}{|X|} \frac{\delta}{2l^2}}_{\leq 1} \mathbb{P}(X_l = V)$$

$$\leq \sum_{l=1}^{\infty} \frac{\delta}{2l^2} \leq \delta$$

$x_*$  is kicked out at round  $l$  if

$$\max_{x \in \mathcal{X}_l} \langle x - x_*, \hat{\theta}_l \rangle > 2\varepsilon_l$$

For any  $x$ , on the good event

$$\langle x - x_*, \hat{\theta}_l \rangle = \langle x - x_*, \hat{\theta}_l - \theta_* \rangle + \langle x - x_*, \theta_* \rangle$$

$$= \underbrace{\langle x, \hat{\theta}_l - \theta_* \rangle}_{\leq \varepsilon_l} - \underbrace{\langle x_*, \hat{\theta}_l - \theta_* \rangle}_{\leq \varepsilon_l} - \underbrace{\Delta_x}_{> 0}$$

$$\leq 2\varepsilon_l.$$

Suppose some  $x \in \mathcal{X}_l$  has  $\Delta_x > 4\varepsilon_l$ ,

then

$$\max_{x' \in \mathcal{X}_l} \langle x' - x, \hat{\theta}_l \rangle \geq \langle x_* - x, \hat{\theta}_l \rangle$$

$$= \langle x_*, \hat{\theta}_l - \theta_* \rangle - \langle x, \hat{\theta}_l - \theta_* \rangle + \underbrace{\Delta_x}_{> 4\varepsilon_l}$$

$$> -2\varepsilon_l + 4\varepsilon_l$$

$$> 2\varepsilon_l \Rightarrow x \text{ is kicked out.}$$

$$\Delta_x > 4\varepsilon_\ell \Rightarrow x \notin \mathcal{X}_{\ell+1} \Rightarrow \Delta_x \leq 4\varepsilon_\ell = 8\varepsilon_{\ell+1}$$

$$\Rightarrow \max_{x \in \mathcal{X}_\ell} \Delta_x \leq 8\varepsilon_\ell \quad \forall \ell.$$

$$R_T = \sum_{x \in \mathcal{X}} \Delta_x \mathbb{E}[T_x]. \quad \text{Fix any } \nu \geq 0$$

$$\sum_{x \in \mathcal{X}} \Delta_x T_x = \sum_{x \in \mathcal{X}: \Delta_x \leq \nu} \Delta_x T_x + \sum_{x \in \mathcal{X}: \Delta_x > \nu} \Delta_x T_x$$

$$\leq \nu T + \sum_{\ell=1}^{\infty} \sum_{x \in \mathcal{X}_\ell: \Delta_x > \nu} \underbrace{\Delta_x}_{\leq 8\varepsilon_\ell} \underbrace{\lceil 3\varepsilon_\ell \hat{\lambda}_{\ell,x} \rceil}_{\leq 3\varepsilon_\ell \hat{\lambda}_{\ell,x} + \mathbb{1}\{\hat{\lambda}_{\ell,x} > 0\}}$$

$a \vee b$   
 $= \max\{a, b\}$

$$\leq \nu T + \sum_{\ell=1}^{\infty} 8\varepsilon_\ell \sum_{x \in \mathcal{X}_\ell: \Delta_x > \nu} (3\varepsilon_\ell \hat{\lambda}_{\ell,x} + \mathbb{1}\{\hat{\lambda}_{\ell,x} > 0\})$$

$\Delta = \min_{x \neq x_*} \langle x_* - x, \theta_* \rangle$

$$\leq \nu T + \sum_{\ell=1}^{\lceil \log_2(8\nu \vee \hat{\Delta}) \rceil} 8\varepsilon_\ell \left( 3\varepsilon_\ell + \frac{(d+1)d}{2} \right)$$

$$\underbrace{\qquad}_{8\varepsilon_\ell^{-2} d \log\left(\frac{4\ell^2 |\mathcal{X}|}{\delta}\right)}$$

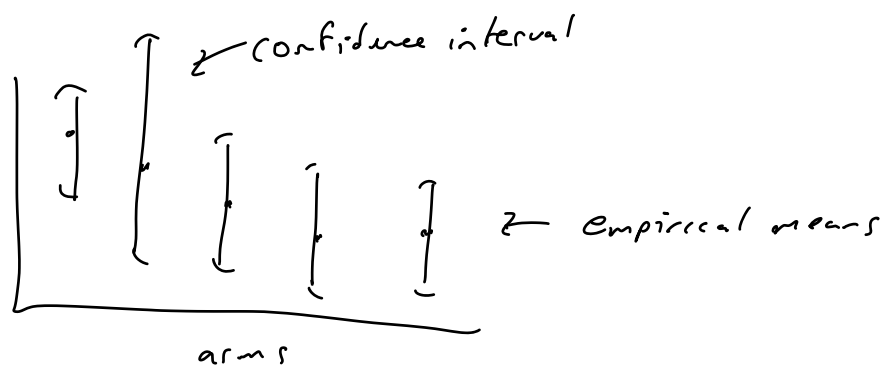
$$\leq \nu T + cd \log\left(\frac{\log_2(8\nu \vee \hat{\Delta}) |\mathcal{X}|}{\delta}\right) \sum_{\ell=1}^{\lceil \log_2(8\nu \vee \hat{\Delta}) \rceil} \underbrace{\varepsilon_\ell^{-1}}_{\frac{1}{2^\ell}} + \sum_{\ell=1}^{\infty} \frac{(d+1)d}{2} \frac{1}{2^\ell}$$

$$\leq \nu T + c' d \log\left(\frac{\log_2(8\nu \vee \hat{\Delta}) |\mathcal{X}|}{\delta}\right) (\nu \vee \Delta)^{-1}$$

To get expected regret choose  $\delta = \frac{1}{T}$

$$\text{If } \nu = 0 \text{ then } R_T \leq \frac{c'd \log(T \log(\delta') |X|)}{\Delta}$$

$\min \nu, \max \Delta \rightarrow R_T \leq c'' \sqrt{dT \log(|X|T)}$



Suppose we've pulled some arms to get  $\{(x_t, y_t)\}_{t \leq S}$  to construct  $\hat{\theta}_S$ . Consider constructing a set

$$C_S = \{\theta \in \mathbb{R}^d : \theta \text{ is not ruled out as } \theta_* \text{ w.p. } \geq 1 - \delta\}$$

UCB Alg.

for  $t=1, 2, \dots$

$$x_t = \operatorname{argmax}_{x \in X} \max_{\theta \in C_{t-1}} \langle x, \theta \rangle$$

# Thompson Sampling

Input prior distribution  $P_0$  defined over  $\mathbb{R}^d$

for  $t = 0, 1, \dots$

$$\tilde{\theta}_t \sim P_t$$

$$x_t = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \langle x, \tilde{\theta}_t \rangle$$

$$\text{Observe } y_t = \langle x_t, \theta^* \rangle + \varepsilon_t$$

Update posterior dist  $P_{t+1}$

$\varepsilon_t \sim \mathcal{D}^2$ -sub-Gaussian

In practice: for some  $\lambda > 0$

$$P_t = \mathcal{N} \left( \left( \sum_{s=1}^t x_s x_s^\top + \lambda I \right)^{-1} \sum_{s=1}^t x_s y_s, \sigma^2 \left( \sum_{s=1}^t x_s x_s^\top + \lambda I \right)^{-1} \right)$$