Fix $x_1, \ldots, x_n \in \mathbb{R}^d$ (known) and $\theta_* \in \mathbb{R}^d$ (unknown).

$\forall i$  Observe

$$y_i = x_i^T \theta_* + \varepsilon_i, \quad \{\varepsilon_i\} \text{ are IID mean-zero}$$

1-sub-Gaussian R.V.

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

**Least squares estimate:**

$$\hat{\theta} = \arg\min_{\theta} \sum_{i=1}^{n} (\langle x_i, \theta \rangle - y_i)^2$$

$$= \left( \sum_{i=1}^{n} x_i x_i^T \right)^{-1} \sum_{i=1}^{n} x_i y_i$$

$$= \left( \sum_{i=1}^{n} x_i x_i^T \right)^{-1} \left( \sum_{i=1}^{n} x_i x_i^T \right) \theta_* + \left( \sum_{i=1}^{n} x_i x_i^T \right)^{-1} \sum_{i=1}^{n} x_i \varepsilon_i$$

$$= \theta_* + \left( \sum_{i=1}^{n} x_i x_i^T \right)^{-1} \sum_{i=1}^{n} x_i \varepsilon_i$$
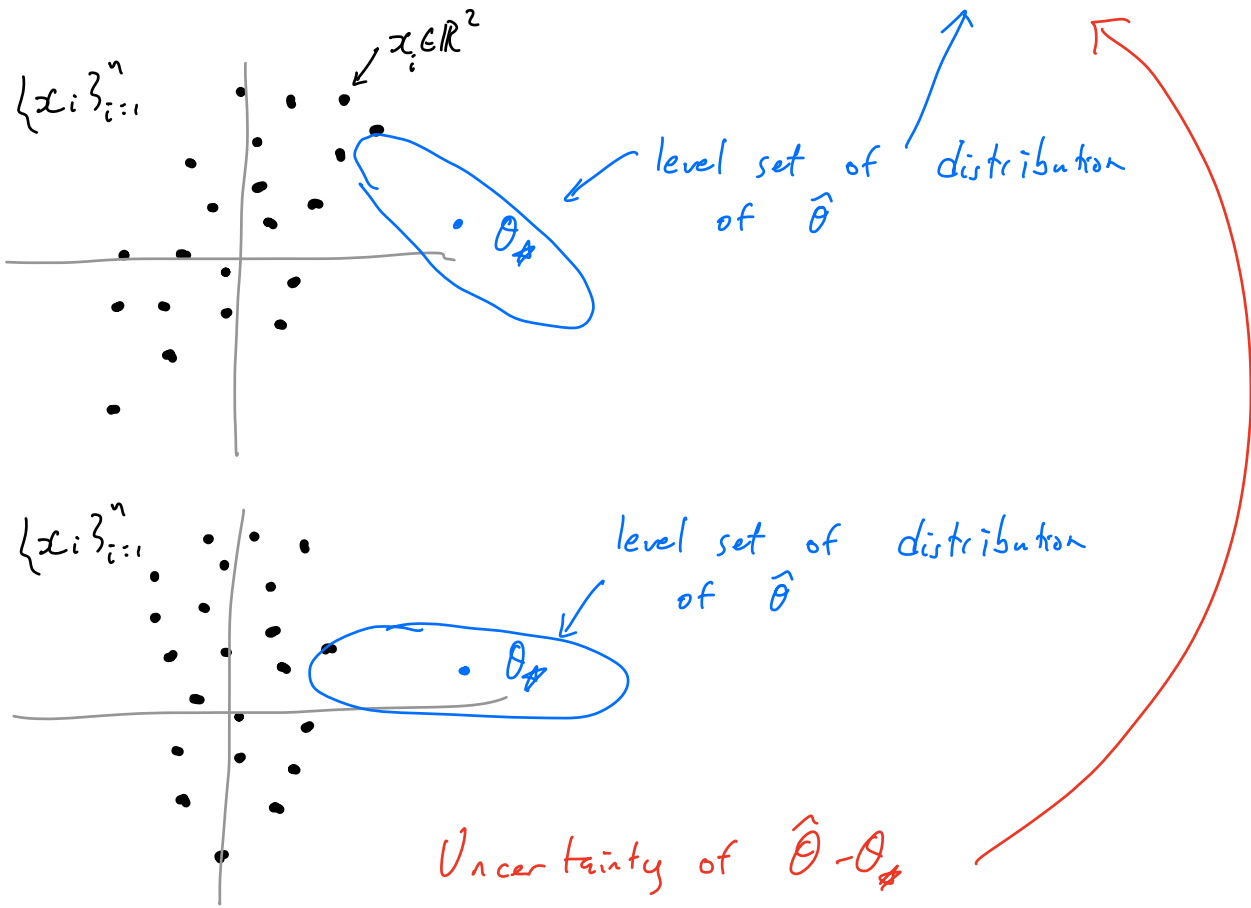
$$\mathbb{E}[\hat{\theta}] = \theta_* \qquad \text{Cov}(\hat{\theta}) = \left( \sum_{i=1}^{n} x_i x_i^T \right)^{-1}$$

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

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

Proof: Show $\langle z, \hat{\theta} - \theta_* \rangle$ is $\|z\|^2_{\left( \sum x_i x_i^T \right)^{-1}}$-sub-Gaussian. Apply Chernoff bound.

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

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

$x_i \in \mathbb{R}^2$

level set of distribution of $\hat{\theta}$

$\theta_*$

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

level set of distribution of $\hat{\theta}$

$\theta_*$

Uncertainty of $\hat{\theta} - \theta_*$

depends only on $\{x_i\}_{i=1}^{n}$ and $\underline{not}$ $\theta_*$ 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 $\hat{\theta}$ w/ $\mathbb{E}[\hat{\theta}] = \theta_*$   $\text{Cov}(\hat{\theta}) = \left(\sum_i x_i x_i^T\right)^{-1}$

Given some pool of points $\mathcal{X} \subset \mathbb{R}^d$
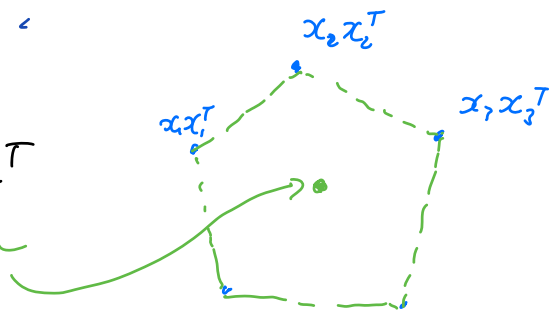
choose $\{x_i\}_{i=1}^n$

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

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

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



Exp. Design Objectives

- A-optimality $\quad f_A(\lambda) = Tr\left( A(\lambda)^{-1} \right)$

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

"minimize average error over all directions"

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

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

"minimize worst-case error direction"

- G - optimality $\quad f_G(\lambda) = \max\limits_{x \in \mathcal{X}} x^T A(\lambda)^{-1} x$

$$\max\limits_{x \in \mathcal{X}} \mathbb{E}\left[\langle x, \hat{\theta} - \theta_* \rangle^2\right] = \max\limits_{x \in \mathcal{X}} x^T \left(\sum_{c=1}^{n} x_i x_i^T\right)^{-1} x$$

$$= \max\limits_{x \in \mathcal{X}} \|x\|^2_{\left(\sum_i x_i x_i^T\right)^{-1}}$$

- D - optimality $\quad f_D(\lambda) = -\log\left(|A(\lambda)|\right)$

$$= \log\left(|A(\lambda)^{-1}|\right)$$

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

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

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

**Lemma** (Kiefer - Wolfowitz 1960)

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

exists a $\lambda^* \in \Delta_{\mathcal{X}}$ such that

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

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

- $f_G(\lambda^*) = d$      Caratheodory Theorem.

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

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

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

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

$s$-sparse if

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

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

$$= d$$

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

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

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

$\lceil 3 \lambda_x^* \rceil$ times, then w.p. $\geq 1 - \delta$

$$\left| \langle x, \hat{\theta} - \theta_* \rangle \right| \leq \sqrt{\frac{2d \log(2|X|/\delta)}{3}} \qquad \forall x \in X$$

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

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

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

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

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

$$= \sum_{x \in X} \lceil 3 \lambda_x^* \rceil x x^T$$

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

$$\underbrace{\phantom{\max_{x' \in X} \|x\|_{\left(\sum_{x \in X} \lambda_x^* x x^T\right)^{-1}}}}_{= \sqrt{d}}$$

$$\geq \sum_{x \in X} 3 \lambda_x^* x x^T$$

$A \succeq B$ if

$A - B \succeq 0$ (PSD)

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

$$\|x\|_{\left(3\sum_{x'\in\mathcal{X}}\lambda_x x'\dot{x}^T\right)^{-1}} = \sqrt{x^T\left(3\sum_{x'\in\mathcal{X}}\lambda_{x'} x'x'^T\right)^{-1}x}$$

$$= \frac{1}{\sqrt{3}} \cdot \|x\|_{A(\lambda^*)^{-1}}$$

**Input**: Finite set $\mathcal{X} \subset \mathbb{R}^d$, confidence level $\delta \in (0,1)$.
Let $\widehat{\mathcal{X}}_1 \leftarrow \mathcal{X}, \ell \leftarrow 1$
**while** $|\mathcal{X}_\ell| > 1$ **do**
  Let $\widehat{\lambda}_\ell \in \triangle_{\mathcal{X}_\ell}$ 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 $x \in \mathcal{X}$ exactly $\lceil \widehat{\lambda}_{\ell,x} \tau_\ell \rceil$ times and construct the least squares estimator $\widehat{\theta}_\ell$ using only the observations of this round
  $\mathcal{X}_{\ell+1} \leftarrow \mathcal{X}_\ell \setminus \{x \in \mathcal{X}_\ell : \max_{x' \in \mathcal{X}_\ell} \langle x' - x, \widehat{\theta}_\ell \rangle > 2\epsilon_\ell \}$
  $\ell \leftarrow \ell + 1$
**Output**: $\mathcal{X}_\ell$

$$\text{Regret} \quad R_T = \max_{x \in \mathcal{X}} \mathbb{E}\left[ \sum_{t=1}^T \langle x, \theta_* \rangle - \langle x_t, \theta_* \rangle \right]$$

$\uparrow$

arm algorithm played at time $t$

$$= \sum_{x \in \mathcal{X}} \Delta_x \, \mathbb{E}[T_x]$$

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

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

Assume $\max_{x \in \mathcal{X}} |\langle x, \theta_* \rangle| \leq 1$. Let $x^* = \underset{x \in \mathcal{X}}{\text{argmax}} \langle x, \theta_* \rangle$

<u>Lemma</u> 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 $R_T \leq \sqrt{d T \log(|\mathcal{X}| T)}$.

I can "cover" $\mathcal{X}$ with just $O(2^d)$ landmarks

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