

Consider a finite set of vectors  $x_1, \dots, x_n \in \mathbb{R}^d$

For all  $i \in [n]$  we observe  $y_i = \langle x_i, \theta^* \rangle + \varepsilon_i$

for  $\varepsilon_i$  <sup>Independent</sup> mean-0, sub-Gaussian,  $\mathbb{E}[\varepsilon_i] = 0$ ,  $\mathbb{E}[\exp(\lambda \varepsilon_i)] \leq \exp(\lambda^2/2)$

Ex.  $y_i \in \{-1, 1\}$  denoting "click/no-click"

$\varepsilon_i \in [-\langle x_i, \theta^* \rangle, 1 - \langle x_i, \theta^* \rangle]$  (bounded)

$\Rightarrow \mathbb{E}[\exp(\lambda \varepsilon_i)] \leq \exp(\lambda^2/2)$

Ex.  $y_i \sim \mathcal{N}(\langle \theta^*, x_i \rangle, 1)$  denoting, "linger time"

Given  $\{(y_i, x_i)\}_{i=1}^n$ , estimate  $\theta^*$ .

$$y_i = \langle x_i, \theta^* \rangle + \varepsilon_i$$

Most natural estimator is LS

(also MLE for Gaussian noise)

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

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

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

$$= (X^T X)^{-1} X^T (X \theta^* + \varepsilon)$$

$$= \theta^* + (X^T X)^{-1} X^T \varepsilon$$

$$X \in \mathbb{R}^{n \times d} \quad y \in \mathbb{R}^n$$

$$X = \begin{pmatrix} \text{---} x_1^T \text{---} \\ \vdots \\ \text{---} x_n^T \text{---} \end{pmatrix}$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

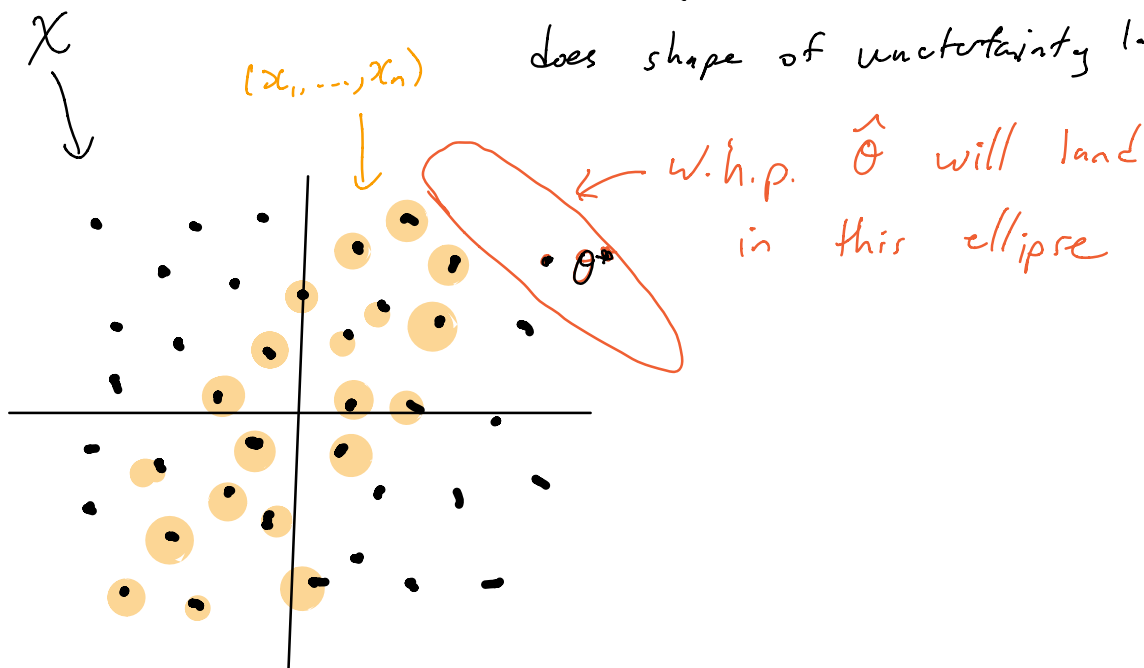
$$y = X \theta^* + \varepsilon$$

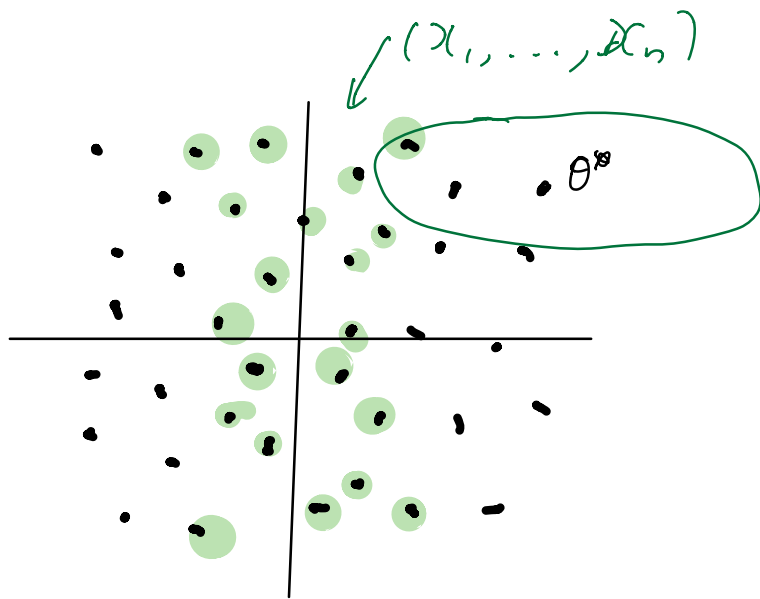
$$\begin{aligned}
\mathbb{E}[(\hat{\theta} - \theta^*)(\hat{\theta} - \theta^*)^T] &= \mathbb{E}[(X^T X)^{-1} X^T \varepsilon \varepsilon^T X (X^T X)^{-1}] \\
&= (X^T X)^{-1} X^T \mathbb{E}[\varepsilon \varepsilon^T] X (X^T X)^{-1} \\
&\leq (X^T X)^{-1} X^T X (X^T X)^{-1} \leftarrow \mathbb{E}[\varepsilon_i^2] \leq 1 \\
&= (X^T X)^{-1} \quad \mathbb{E}[\varepsilon_i \varepsilon_j] = 0
\end{aligned}$$

Ex. If  $\varepsilon \sim \mathcal{N}(0, 1)$  then

$\hat{\theta}$  has mean  $\theta^*$ , covariance  $(X^T X)^{-1}$   
 $\hat{\theta} - \theta^* \sim \mathcal{N}(0, (X^T X)^{-1})$

Now suppose we have a pool of points  $\mathcal{X}$  and we choose  $(x_1, \dots, x_n)$  from it. How does shape of uncertainty look?





For any PSD matrix  $A$  define  $\|x\|_A^2 = x^T A x$ .

$A$  is PSD if  $x^T A x \geq 0$  for all  $x \in \mathbb{R}^d$ .

Proposition | Fix  $x_1, \dots, x_n$ . Let  $y_i = \langle x_i, \theta^* \rangle + \varepsilon_i$   $\theta^* \in \mathbb{R}^d$  where  $\mathbb{E}[\varepsilon_i] = 0$  and  $\mathbb{E}(\exp(\lambda \varepsilon_i)) \leq e^{\lambda^2/2}$  and  $\varepsilon_i \perp \varepsilon_j$ . Then for any  $z \in \mathbb{R}^d$  w.p.  $\geq 1 - \delta$ ,  $\hat{\theta} = \arg \min_{\theta} \sum_i (y_i - x_i^T \theta)^2$  satisfies

$$\langle z, \hat{\theta} - \theta^* \rangle \leq \sqrt{2 \|z\|_{(X^T X)^{-1}}^2 \log(1/\delta)}.$$

Which follows from  $\mathbb{P}(\langle z, \hat{\theta} - \theta^* \rangle \geq t) \leq \exp(-t^2 / 2 \|z\|_{(X^T X)^{-1}}^2)$ .

Proof)  $\hat{\theta} = \theta_* + (X^T X)^{-1} X^T \varepsilon$   $w_i = X (X^T X)^{-1} z$ .

$$\Rightarrow z^T (\hat{\theta} - \theta^*) = z^T (X^T X)^{-1} X^T \varepsilon =: w^T \varepsilon = \sum_{i=1}^n w_i \varepsilon_i$$

$$\begin{aligned}
\mathbb{E}[\exp(\lambda \sum_{i=1}^n \omega_i \varepsilon_i)] &= \mathbb{E}\left[\prod_{i=1}^n \exp(\lambda \omega_i \varepsilon_i)\right] \\
&= \prod_{i=1}^n \mathbb{E}[\exp(\lambda \omega_i \varepsilon_i)] && \text{(Independence of } \varepsilon_i) \\
&\leq \prod_{i=1}^n \exp(\lambda^2 \omega_i^2 / 2) && \text{(sub-Gaussian)} \\
&= \exp\left(\sum_{i=1}^n \lambda^2 \omega_i^2 / 2\right) \\
&= \exp(\lambda^2 \|\omega\|_2^2 / 2)
\end{aligned}$$

Through Chernoff bound technique

$$\mathbb{P}(\omega^T \varepsilon \geq t) \leq \exp\left(-\frac{t^2}{2\|\omega\|_2^2}\right)$$

$$\begin{aligned}
\|\omega\|_2^2 &= \|X(X^T X)^{-1} z\|_2^2 \\
&= z^T (X^T X)^{-1} \cancel{X^T X} (X^T X)^{-1} z \\
&= z^T (X^T X)^{-1} z = \|z\|_{(X^T X)^{-1}}^2
\end{aligned}$$

$$\begin{aligned}
\mathbb{P}(\langle \hat{\theta} - \theta^*, z \rangle \geq t) &= \mathbb{P}(\langle \omega, \varepsilon \rangle \geq t) \\
&\leq \exp\left(-\frac{t^2}{2\|\omega\|_2^2}\right) \\
&= \exp\left(-\frac{t^2}{2\|z\|_{(X^T X)^{-1}}^2}\right)
\end{aligned}$$

□

So given cov of  $\hat{\theta} - \theta^*$  is  $(X^T X)^{-1}$

$$\text{and } \langle z, \hat{\theta} - \theta^* \rangle \leq \sqrt{z^T (X^T X)^{-1} z} \log(1/\delta).$$

Clearly,  $(X^T X)^{-1}$  is key quantity.

For any  $x_1, \dots, x_n \in \mathcal{X} \exists \lambda \in \Delta_{\mathcal{X}} := \{ \nu \in \mathbb{R}_+^{|\mathcal{X}|} : \sum_{x \in \mathcal{X}} \nu_x = 1 \}$

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

↑  
pool of points

$$A(\lambda) = \sum_{x \in \mathcal{X}} \lambda_x x x^T \leftarrow \text{Design matrix.}$$

Idea: instead of choosing  $(x_1, \dots, x_n)$ , simply design  $\lambda \in \Delta_{\mathcal{X}}$  ignoring integer effects.

### Design Criterion

A-optimal:  $f_A(\lambda) = \text{Tr}(A(\lambda)^{-1})$  to minimize  $\mathbb{E}[\|\hat{\theta} - \theta^*\|_2^2]$

$$\mathbb{E}[\|\hat{\theta} - \theta^*\|_2^2] = \mathbb{E}[\|(X^T X)^{-1} X^T \epsilon\|_2^2] \leq \text{Tr}((X^T X)^{-1}) \quad (\text{equality if } \mathbb{E}[\epsilon_i^2] = 1)$$

$$\approx \frac{1}{n} \text{Tr}(A(\lambda)^{-1})$$

E-optimal:  $f_E(\lambda) = \lambda_{\max}(A(\lambda)^{-1})$  to min.  $\sup_{\|u\| \leq 1} E[\langle u, \hat{\theta} - \theta_* \rangle^2]$

D-optimal:  $g_D(\lambda) = \log |A(\lambda)|$ . Maximizing this quantity is equivalent to minimizing the entropy of  $\hat{\theta}$  if  $\varepsilon_i \sim \mathcal{N}(0, 1)$ .

G-optimal:  $f_G(\lambda) = \max_{x \in \mathcal{X}} \|x\|_{A(\lambda)^{-1}}^2$  to min.  $\max_{x \in \mathcal{X}} E[\langle x, \hat{\theta} - \theta_* \rangle^2]$

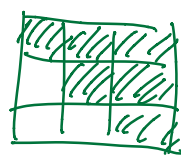
### Theorem / Kiefer-Wolfowitz (1960)

For finite set  $\mathcal{X}$  w/  $\dim(\text{span}(\mathcal{X})) = d$  there exists  $\lambda^* \in \Delta_{\mathcal{X}}$  such that:

- $\max_{\lambda} g_D(\lambda) = g_D(\lambda^*)$
- $\min_{\lambda} f_G(\lambda) = f_G(\lambda^*)$
- $f_G(\lambda^*) = g_D(\lambda^*) = d$ .
- $\text{support}(\lambda^*) \leq \frac{(d+1)d}{2}$

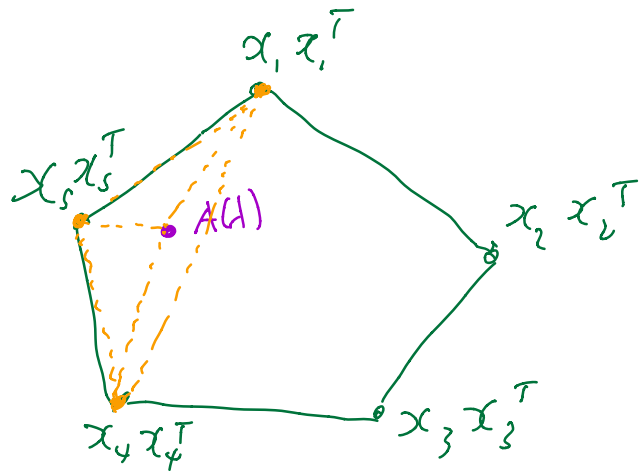
Fix any  $\lambda \in \Delta_{\mathcal{X}}$  and consider

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



degrees of freedom  $\leq \frac{(d+1)d}{2}$

Consider each  $x_i x_i^T$  as a point  
in  $\mathbb{R}^{(d+1)d/2}$



In picture  $\dim p=2$

Carathéodory says any pt w/in convex  
hull of  $n$  points in  $p$ -dimensions  
can be represented by a convex combination  
of just  $p+1$  points.

Proposition Fix  $\mathfrak{S} \in \mathbb{N}$  and let  $\lambda^*$  be the  $G$ -optimal soln w/  $\text{support}(\lambda^*) \leq \frac{d(d+1)}{2}$ . Pull arm  $x \in \mathcal{X}$  exactly  $\lceil \mathfrak{S} \lambda_x^* \rceil$  times and compute  $\hat{\theta} = (X^T X)^{-1} X^T y$  where  $X$  are stacked vectors and  $y$  are observations. Then w.p.  $\geq 1 - \delta$  for all  $x \in \mathcal{X}$  we have

$$|\langle \hat{\theta} - \theta^*, x \rangle| \leq \sqrt{2 \|x\|_{(X^T X)^{-1}}^2 \log(2|\mathcal{X}|/\delta)}$$

$$= \sqrt{x^T \left( \sum_{x' \in \mathcal{X}} \lceil \mathfrak{S} \lambda_{x'}^* \rceil x' x'^T \right)^{-1} x \cdot 2 \log(2|\mathcal{X}|/\delta)}$$

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

$$= \|x\|_{A(\lambda^*)^{-1}} \sqrt{\frac{2 \log(2|\mathcal{X}|/\delta)}{\mathfrak{S}}}$$

$$\leq \sqrt{\frac{2d \log(2|\mathcal{X}|/\delta)}{\mathfrak{S}}}$$

Note: total pulls  $\leq \mathfrak{S} + \text{support}(\lambda^*) \leq \mathfrak{S} + \frac{(d+1)d}{2}$ .