

# Minimax experimental design

Bridging the gap between statistical and worst-case approaches to least squares regression

Michał Dereziński

Kenneth L. Clarkson

Michael W. Mahoney

Manfred K. Warmuth

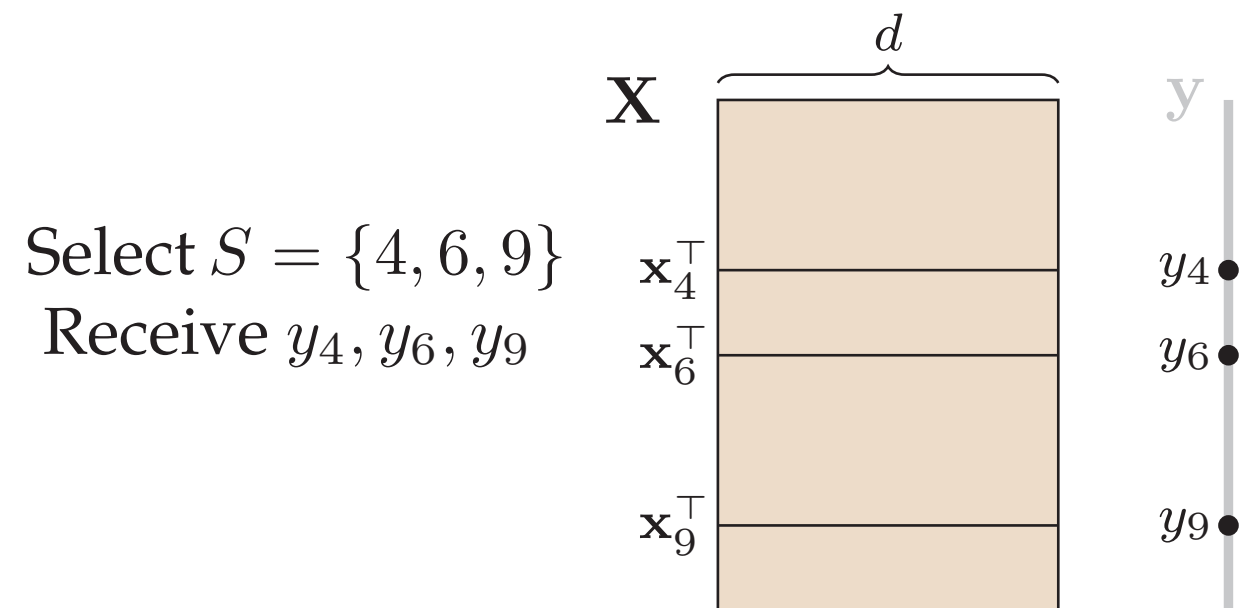


## Classical experimental design

Consider  $n$  vectors:  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ .  
Each  $\mathbf{x}_i$  has a hidden random response  $y_i$ :

$$y_i = \mathbf{x}_i^\top \mathbf{w}^* + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma^2)$$

**Task:** Select  $k \ll n$  responses to query



**Goal:** Find the *best unbiased estimator* of  $\mathbf{w}^*$

## A-optimal design

Minimize the *mean squared error*:

$$\min_{\hat{\mathbf{w}}} \underbrace{\mathbb{E}_{\hat{\mathbf{w}}} [\|\hat{\mathbf{w}} - \mathbf{w}^*\|^2]}_{\text{MSE}[\hat{\mathbf{w}}]} \quad \text{s.t.} \quad \mathbb{E}[\hat{\mathbf{w}}] = \mathbf{w}^*$$

**Classical result:** given  $\{y_i : i \in S\}$ ,  
the optimum is *least squares*,  $\hat{\mathbf{w}} = \mathbf{X}_S^\dagger \mathbf{y}_S$

$$\begin{aligned} \text{MSE}[\mathbf{X}_S^\dagger \mathbf{y}_S] &= \text{tr}(\text{Var}[\mathbf{X}_S^\dagger \mathbf{y}_S]) \\ &= \sigma^2 \text{tr}((\mathbf{X}_S^\top \mathbf{X}_S)^{-1}) \end{aligned}$$

A-optimal design:  $\min_{S: |S| \leq k} \text{tr}((\mathbf{X}_S^\top \mathbf{X}_S)^{-1})$

**Alternative:** Mean squared prediction error

$$\text{MSPE}[\hat{\mathbf{w}}] = \mathbb{E}_{\hat{\mathbf{w}}} [\|\mathbf{X}(\hat{\mathbf{w}} - \mathbf{w}^*)\|^2]$$

V-optimal design:  $\min_{S: |S| \leq k} \text{tr}(\mathbf{X}(\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}^\top) \propto \text{MSPE}[\mathbf{X}_S^\dagger \mathbf{y}_S]$

All of our results extend to V-optimal design

## Prior: A simple guarantee

**Theorem** (from [?]) For any  $\mathbf{X}$  and  $k \geq d$  there is  $S$  of size  $k$  s.t.:

$$\text{tr}((\mathbf{X}_S^\top \mathbf{X}_S)^{-1}) \leq \frac{n-d+1}{k-d+1} \underbrace{\text{tr}((\mathbf{X}^\top \mathbf{X})^{-1})}_{\text{(denoted } \phi)}$$

Let  $\boldsymbol{\xi}^\top = [\xi_1, \dots, \xi_n]$  be the noise vector

**Corollary** If  $\text{Var}[\boldsymbol{\xi}] = \sigma^2 \mathbf{I}$  and  $\mathbb{E}[\boldsymbol{\xi}] = \mathbf{0}$  then

$$\begin{aligned} \text{tr}(\text{Var}[\mathbf{X}_S^\dagger \mathbf{y}_S]) &\leq \sigma^2 \frac{n-d+1}{k-d+1} \phi \\ &\leq \underbrace{\frac{\phi}{k-d+1}}_{\epsilon} \cdot \underbrace{\text{tr}(\text{Var}[\boldsymbol{\xi}])}_{n\sigma^2} \end{aligned}$$

For  $k = d + \phi/\epsilon$  there is  $S$  of size  $k$  s.t.  
 $\text{MSE}[\mathbf{X}_S^\dagger \mathbf{y}_S] \leq \epsilon \cdot \text{tr}(\text{Var}[\boldsymbol{\xi}])$

**Question:** What if we drop the assumptions on the noise vector  $\boldsymbol{\xi}$ ?

## Arbitrary response model

$\mathbf{y} \in \mathcal{F}_n$  - random vector in  $\mathbb{R}^n$  s.t.  $\mathbb{E}[\|\mathbf{y}\|^2] < \infty$

$$\mathbf{w}_{\mathbf{y}|\mathbf{X}}^* \stackrel{\text{def}}{=} \underset{\mathbf{w}}{\text{argmin}} \mathbb{E}_{\mathbf{y}} [\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2]$$

$$\boldsymbol{\xi}_{\mathbf{y}|\mathbf{X}} \stackrel{\text{def}}{=} \mathbf{y} - \mathbf{X}\mathbf{w}_{\mathbf{y}|\mathbf{X}}^*$$

Two special cases:

1. Statistical regression:  $\mathbb{E}[\boldsymbol{\xi}_{\mathbf{y}|\mathbf{X}}] = \mathbf{0}$

2. Worst-case regression:  $\text{Var}[\boldsymbol{\xi}_{\mathbf{y}|\mathbf{X}}] = \mathbf{0}$

## Random experimental designs

**Definition 1** We define a random experimental design  $(S, \hat{\mathbf{w}})$  of size  $k$  as:

1. a random  $S \subseteq \{1..n\}$  s.t.  $|S| \leq k$
2. a jointly random function  $\hat{\mathbf{w}} : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^d$

$\mathcal{W}_k(\mathbf{X})$  - family of *unbiased* random experimental designs  $(S, \hat{\mathbf{w}})$  of size  $k$ :

$$\mathbb{E}_{S, \hat{\mathbf{w}}, \mathbf{y}} [\hat{\mathbf{w}}(\mathbf{y}_S)] = \mathbf{w}_{\mathbf{y}|\mathbf{X}}^* \quad \text{for all } \mathbf{y} \in \mathcal{F}_n$$

## Main result: experimental design with arbitrary responses

**Theorem 1** There is a random experimental design of size

$$k = O(d \log n + \phi/\epsilon), \quad \text{where } \phi = \text{tr}((\mathbf{X}^\top \mathbf{X})^{-1}),$$

such that for any random response vector  $\mathbf{y} \in \mathcal{F}_n$  we have

$$\text{(unbiasedness)} \quad \mathbb{E}[\hat{\mathbf{w}}(\mathbf{y}_S)] = \mathbf{w}_{\mathbf{y}|\mathbf{X}}^*,$$

$$\text{MSE}[\hat{\mathbf{w}}(\mathbf{y}_S)] - \text{MSE}[\mathbf{X}_S^\dagger \mathbf{y}_S] \leq \epsilon \cdot \mathbb{E}[\|\boldsymbol{\xi}_{\mathbf{y}|\mathbf{X}}\|^2].$$

### Important examples

1. *Statistical regression:*

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \boldsymbol{\xi}, \quad \mathbb{E}[\boldsymbol{\xi}] = \mathbf{0}$$

- *Weighted regression:*  
 $\text{Var}[\boldsymbol{\xi}] = \text{diag}([\sigma_1^2, \dots, \sigma_n^2])$
- *Generalized regression:*  
 $\text{Var}[\boldsymbol{\xi}]$  is arbitrary

2. *Worst-case regression:*

$\mathbf{y}$  is any fixed vector in  $\mathbb{R}^n$

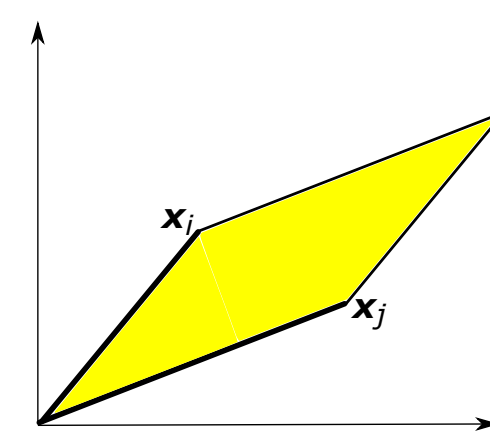
## Volume sampling

**Definition 2** Given a full rank matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  we define volume sampling  $\text{VS}(\mathbf{X})$  as a distribution over sets  $S \subseteq [n]$  of size  $d$ :

$$\Pr(S) = \frac{\det(\mathbf{X}_S)^2}{\det(\mathbf{X}^\top \mathbf{X})}$$

$\Pr(S) \sim$  squared volume of the parallelepiped spanned by  $\{\mathbf{x}_i : i \in S\}$

Computational cost:  $O(nd^2)$



**Theorem** (from [?]) Volume sampling corrects the bias of any i.i.d. sampling  $q = (q_1, \dots, q_n)$ :

$$\text{volume + i.i.d.} \quad \underbrace{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_d}}_{\sim \text{VS}(\mathbf{X})}, \underbrace{\mathbf{x}_{i_{d+1}}, \dots, \mathbf{x}_{i_k}}_{\sim q^{k-d}}$$

$$\mathbb{E} \left[ \underset{\mathbf{w}}{\text{argmin}} \sum_{t=1}^k \frac{1}{q_{i_t}} (\mathbf{x}_{i_t}^\top \mathbf{w} - y_{i_t})^2 \right] = \mathbf{w}_{\mathbf{y}|\mathbf{X}}^*$$

## Leverage and inverse scores

$\mathbf{z}_i$  -  $i$ th column of pseudoinverse  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

1. *Leverage score sampling:*

$$\Pr(i) = p_i^{\text{lev}} \stackrel{\text{def}}{=} \frac{1}{d} \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i = \frac{1}{d} \mathbf{x}_i^\top \mathbf{z}_i$$

used to obtain a subspace embedding

2. *Inverse score sampling (new):*

$$\Pr(i) = p_i^{\text{inv}} \stackrel{\text{def}}{=} \frac{1}{\phi} \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-2} \mathbf{x}_i = \frac{1}{\phi} \|\mathbf{z}_i\|^2$$

used to get the  $\phi/\epsilon$  rate for the MSE

Controlling the least squares estimator:

$$\begin{aligned} \hat{\mathbf{w}} &= \underset{\mathbf{w}}{\text{argmin}} \sum_{t=1}^k \frac{1}{q_{i_t}} (\mathbf{x}_{i_t}^\top \mathbf{w} - y_{i_t})^2 \\ &= \underbrace{\left( \sum_{t=1}^k \frac{1}{q_{i_t}} \mathbf{z}_{i_t} \mathbf{x}_{i_t}^\top \right)^{-1}}_{\text{needs leverage scores}} \underbrace{\sum_{t=1}^k \frac{1}{q_{i_t}} \mathbf{z}_{i_t} y_{i_t}}_{\text{needs inverse scores}} \end{aligned}$$

**Strategy:** Use a mixture of both  $p^{\text{lev}}$  and  $p^{\text{inv}}$

## Minimax A-optimal value for experimental design

**Definition 3** Define the minimax A-optimal value for experimental design with unbiased estimators as:

$$R_k^*(\mathbf{X}) \stackrel{\text{def}}{=} \min_{(S, \hat{\mathbf{w}}) \in \mathcal{W}_k(\mathbf{X})} \max_{\mathbf{y} \in \mathcal{F}_n \setminus \text{Sp}(\mathbf{X})} \frac{\text{MSE}[\hat{\mathbf{w}}(\mathbf{y}_S)] - \text{MSE}[\mathbf{X}_S^\dagger \mathbf{y}_S]}{\mathbb{E}[\|\boldsymbol{\xi}_{\mathbf{y}|\mathbf{X}}\|^2]}$$

**Proposition 2** Least squares is the minimum variance unbiased estimator for  $\mathcal{F}_n$ :

$$\begin{aligned} \text{if } \mathbb{E}_{\mathbf{y}, \hat{\mathbf{w}}} [\hat{\mathbf{w}}(\mathbf{y})] &= \mathbf{X}^\dagger \mathbb{E}[\mathbf{y}] \quad \forall \mathbf{y} \in \mathcal{F}_n, \\ \text{then } \text{Var}[\hat{\mathbf{w}}(\mathbf{y})] &\succeq \text{Var}[\mathbf{X}_S^\dagger \mathbf{y}_S] \quad \forall \mathbf{y} \in \mathcal{F}_n. \end{aligned}$$

**Theorem 3** The following statements are true about the minimax A-optimal value:

1. If  $d \leq k \leq n$ , then  $R_k^*(\mathbf{X}) \in [0, \infty)$  for a full rank  $\mathbf{X}$ .
2. If  $k \geq C \cdot d \log n$ , then  $R_k^*(\mathbf{X}) \leq C \cdot \phi/k$  for some  $C$ . (bound from Theorem ??)
3. If  $k^2 < \epsilon n d/3$ , then  $R_k^*(\mathbf{X}) \geq (1-\epsilon) \cdot \phi/k$  for some  $\mathbf{X}$ . (matching lower bound)

### References

- [1] H. Avron, C. Boutsidis. *Faster subset selection for matrices and applications*. JMAA 2013.
- [2] M. Dereziński, M. Warmuth, D. Hsu. *Correcting the bias in least squares regression with volume-rescaled sampling*. AISTATS'19.