

Correcting the bias in least squares regression with volume-rescaled sampling

Michał Dereziński



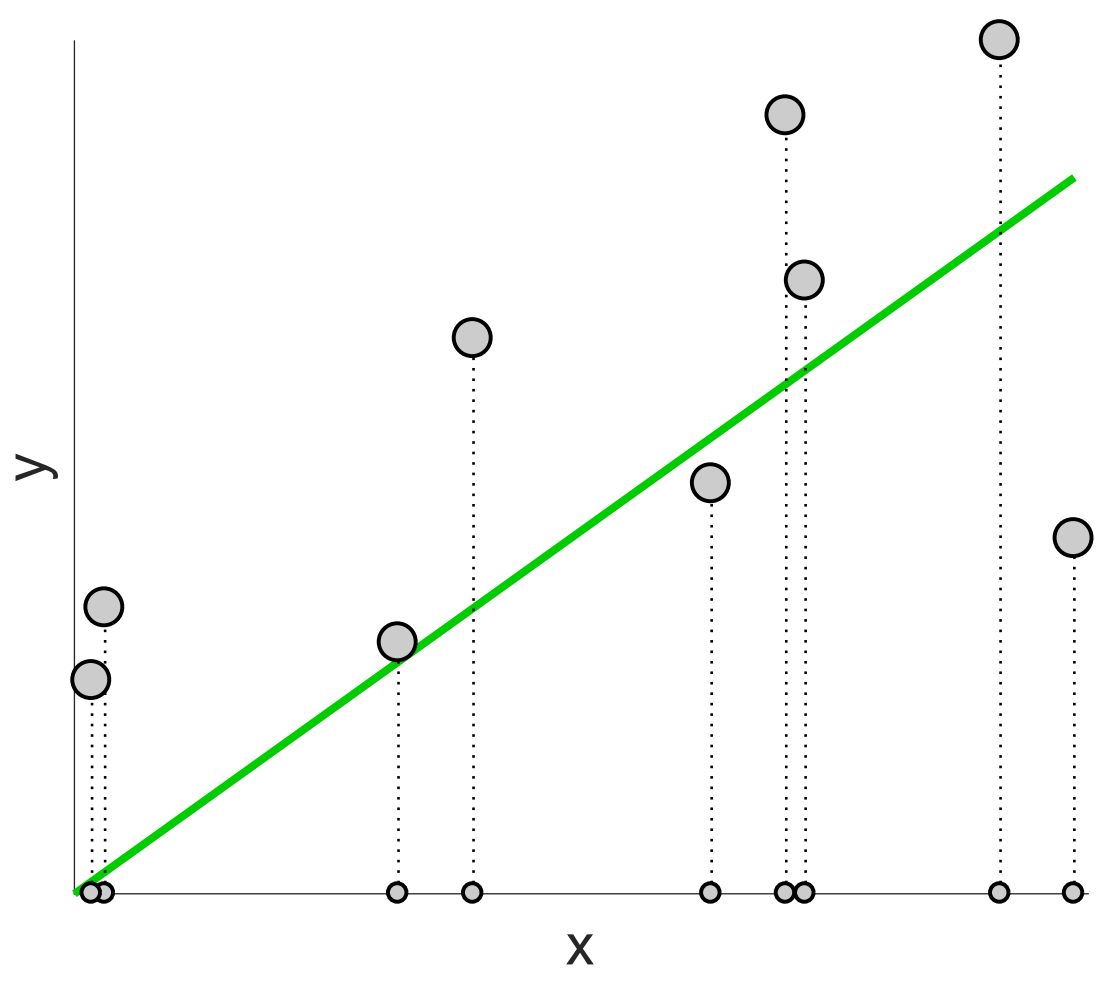
Manfred K. Warmuth



Daniel Hsu



Random design regression



Dataset: $\mathbb{S} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$

Optimum: $\mathbf{w}_D^* = \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{\mathcal{D}} [(\mathbf{x}^\top \mathbf{w} - y)^2]$

Estimator: $\mathbf{w}^*(\mathbb{S}) = \operatorname{argmin}_{\mathbf{w}} \sum_{(\mathbf{x}_i, y_i) \in \mathbb{S}} (\mathbf{x}_i^\top \mathbf{w} - y_i)^2$

Bias of the least squares

$\mathbf{x}_i \sim \mathcal{D}_{\mathcal{X}}$ - random input vector in \mathbb{R}^d

$y_i \sim \mathcal{D}_{\mathcal{Y}|\mathbf{x}=\mathbf{x}_i}$ - random response variable in \mathbb{R}

$$\mathbf{w}^*(\mathbb{S}) = \left(\underbrace{\sum_i \mathbf{x}_i \mathbf{x}_i^\top}_{\text{inverse covariance}} \right)^{-1} \sum_i \mathbf{x}_i y_i$$

Random inverse covariance introduces **bias**:

$$\mathbb{E}[\mathbf{w}^*(\mathbb{S})] \neq \mathbf{w}_D^*$$

Correcting the bias

Our approach: Augment dataset with d points drawn with *volume-rescaled sampling*

Sample d points $\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+d} \sim \det \begin{pmatrix} -\mathbf{x}_{n+1}^\top \\ \vdots \\ -\mathbf{x}_{n+d}^\top \end{pmatrix} \cdot \mathcal{D}_{\mathcal{X}}^d$

Query $y_{n+i} \sim \mathcal{D}_{\mathcal{Y}|\mathbf{x}=\mathbf{x}_{n+i}} \quad \forall i=1..d$

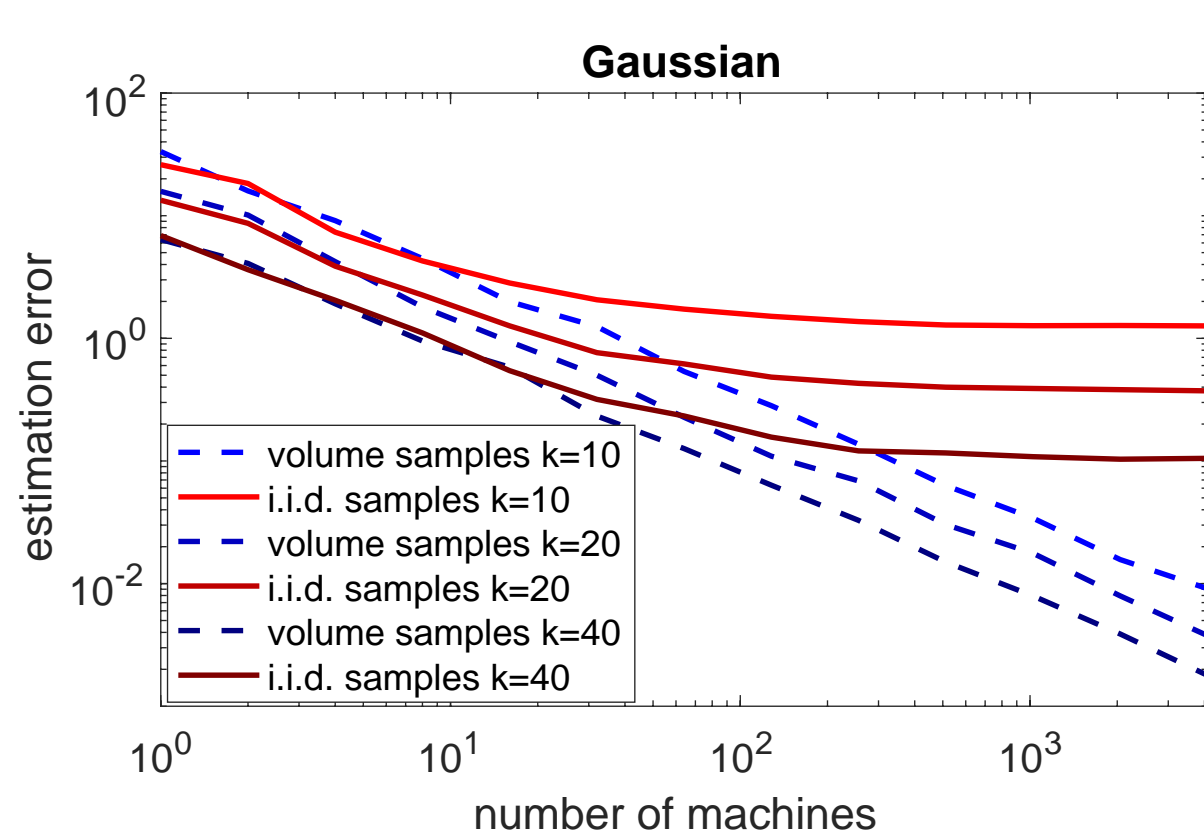
Add $\mathbb{S}_o = (\mathbf{x}_{n+1}, y_{n+1}), \dots, (\mathbf{x}_{n+d}, y_{n+d})$

Theorem $\mathbb{E}[\mathbf{w}^*(\langle \mathbb{S}, \mathbb{S}_o \rangle)] = \mathbf{w}^*$

Example: Averaging estimators

$\mathbf{x}^\top = (x_1, \dots, x_d) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1), \quad y = \xi(\mathbf{x}) + \epsilon,$
 $\xi(\cdot)$ is a non-linear function and $\epsilon \sim \mathcal{N}(0, 1)$.

estimation error: $\left\| \frac{1}{T} \sum_{t=1}^T \mathbf{w}^*(\mathbb{S}_t) - \mathbf{w}_D^* \right\|^2$



Plot uses $d = 5$ and $\xi(\mathbf{x}) = \sum_{i=1}^d x_i + \frac{x_i^3}{3}$

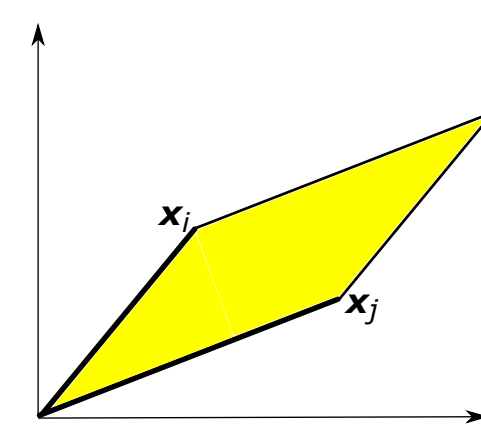
Volume-rescaled sampling

Definition 1 Given $\mathcal{D}_{\mathcal{X}}$ and any $k \geq d$,

$$\text{VS}_{\mathcal{D}_{\mathcal{X}}}^k(\mathbf{x}_1, \dots, \mathbf{x}_k) \propto \det \left(\sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \right) \cdot \prod_{i=1}^k \mathcal{D}_{\mathcal{X}}(\mathbf{x}_i)$$

When $k = d$, then

$$\det \left(\sum_{i=1}^d \mathbf{x}_i \mathbf{x}_i^\top \right) = \det \begin{pmatrix} -\mathbf{x}_1^\top \\ \vdots \\ -\mathbf{x}_d^\top \end{pmatrix}^2 = \text{volume}^2 \dots$$

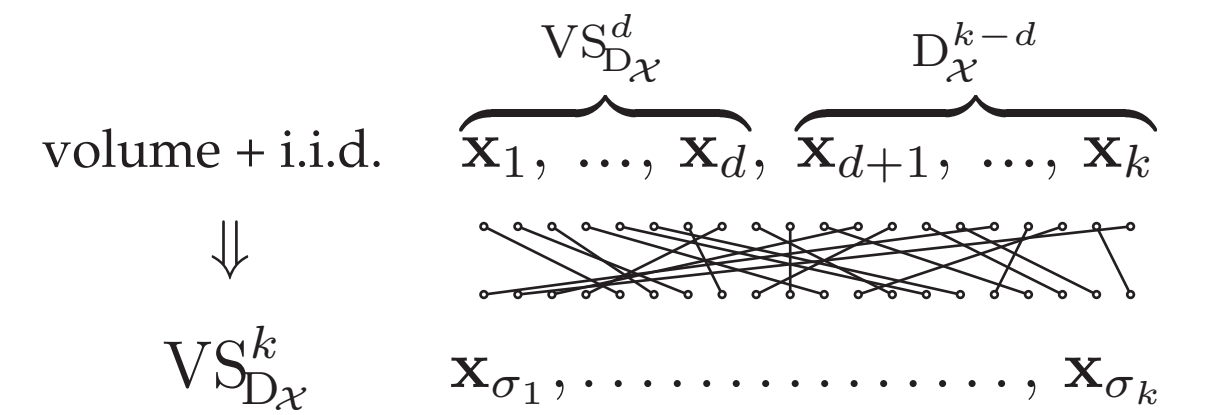


... of the parallelepiped spanned by $\mathbf{x}_1, \dots, \mathbf{x}_d$

Augmentation property

Theorem 1 Let $\mathbb{S}_o \sim \text{VS}_{\mathcal{D}_{\mathcal{X}}}^d, \mathbb{S} \sim \mathcal{D}_{\mathcal{X}}^{k-d}$ and σ be a uniformly random permutation of $\{1..n\}$. Then

$$\sigma(\langle \mathbb{S}_o, \mathbb{S} \rangle) \sim \text{VS}_{\mathcal{D}_{\mathcal{X}}}^k$$

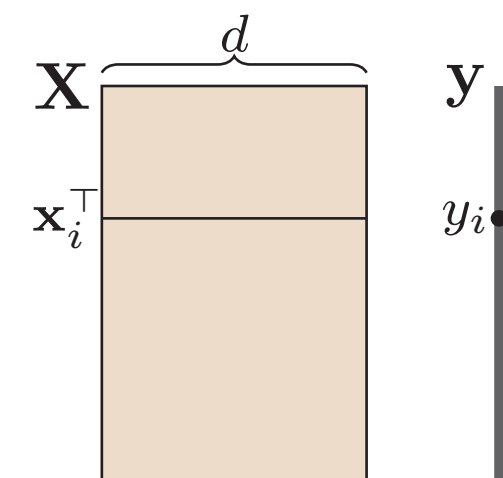


Unbiasedness via augmentation

$(\mathbf{X}, \mathbf{y}) \sim \text{VS}_{\mathcal{D}}^k$:

$\mathbf{x}_1, \dots, \mathbf{x}_k \sim \text{VS}_{\mathcal{D}_{\mathcal{X}}}^k,$

$y_i \sim \mathcal{D}_{\mathcal{Y}|\mathbf{x}=\mathbf{x}_i} \quad \forall i.$



$\mathbb{S} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_k, y_k) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$
 $\mathbb{S}_o = (\tilde{\mathbf{x}}_1, \tilde{y}_1), \dots, (\tilde{\mathbf{x}}_d, \tilde{y}_d) \sim \text{VS}_{\mathcal{D}}^d$

$$\mathbb{E}[\mathbf{w}^*(\langle \mathbb{S}_o, \mathbb{S} \rangle)] = \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^k} [\mathbb{E}_{\mathbb{S}_o \sim \text{VS}_{\mathcal{D}}^d} [\mathbf{w}^*(\langle \mathbb{S}_o, \mathbb{S} \rangle)]]$$

(Theorem 1) $= \mathbb{E}_{\tilde{\mathbb{S}} \sim \text{VS}_{\mathcal{D}}^{k+d}} [\mathbf{w}^*(\tilde{\mathbb{S}})]$

(Theorem 2) $= \mathbf{w}_D^*$

Theorem 2 For any distribution \mathcal{D} ,

$$\mathbb{E}_{\text{VS}_{\mathcal{D}}^k} [(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}] = \left(\mathbb{E}_{\mathcal{D}_{\mathcal{X}}} [\mathbf{x} \mathbf{x}^\top] \right)^{-1} \mathbb{E}_{\mathcal{D}} [\mathbf{x} y]$$

Sampling complexity

Conditioning number of $\mathcal{D}_{\mathcal{X}}$:

$$K_{\mathcal{D}_{\mathcal{X}}} = \sup_{\tilde{\mathbf{x}} \in \text{supp}(\mathcal{D}_{\mathcal{X}})} \tilde{\mathbf{x}}^\top \left(\mathbb{E}_{\mathcal{D}_{\mathcal{X}}} [\mathbf{x} \mathbf{x}^\top] \right)^{-1} \tilde{\mathbf{x}}$$

Theorem 3 The complexity of $\text{VS}_{\mathcal{D}_{\mathcal{X}}}^d$ is:

1. $O(K_{\mathcal{D}_{\mathcal{X}}} d \log d)$ i.i.d. samples from $\mathcal{D}_{\mathcal{X}}$,

2. $O(K_{\mathcal{D}_{\mathcal{X}}} d^3 \log d)$ arithmetic operations

Gaussian algorithm

Theorem 4 Suppose that $\mathcal{D}_{\mathcal{X}} = \mathcal{N}(\mathbf{0}, \Sigma)$.

Let $\mathbf{x}_1, \dots, \mathbf{x}_{2d+2} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_{\mathcal{X}}$ and

$$\tilde{\mathbf{x}}_i \stackrel{\text{def}}{=} \left(\sum_{j=d+1}^{2d+2} \mathbf{x}_j \mathbf{x}_j^\top \right)^{\frac{1}{2}} \left(\sum_{j=1}^d \mathbf{x}_j \mathbf{x}_j^\top \right)^{-\frac{1}{2}} \mathbf{x}_i.$$

Then $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_d \sim \text{VS}_{\mathcal{D}_{\mathcal{X}}}^d$

General volume-rescaled sampling algorithm

Construct an approximate covariance matrix $\hat{\Sigma}$:

$$(1 - \epsilon) \cdot \mathbb{E}_{\mathcal{D}_{\mathcal{X}}} [\mathbf{x} \mathbf{x}^\top] \preceq \hat{\Sigma} \preceq (1 + \epsilon) \cdot \mathbb{E}_{\mathcal{D}_{\mathcal{X}}} [\mathbf{x} \mathbf{x}^\top], \quad \text{where } \epsilon = \frac{1}{\sqrt{2d}}$$

Leverage score sampling w.r.t. matrix $\hat{\Sigma}$:

$$\text{Lev}_{\hat{\Sigma}, \mathcal{X}}(\mathbf{x}) \propto l_{\hat{\Sigma}}(\mathbf{x}) \cdot \mathcal{D}_{\mathcal{X}}(\mathbf{x}), \quad \text{where } l_{\hat{\Sigma}}(\mathbf{x}) = \mathbf{x}^\top \hat{\Sigma}^{-1} \mathbf{x}$$

Key idea: Distortion-free intermediate sampling [DWH18, Der19]

- Sample a larger intermediate sample **here:** d^2 i.i.d. leverage score samples
- Downsample to the target size **here:** d samples

Does not distort the target distribution even though $\hat{\Sigma} \neq \mathbb{E}_{\mathcal{D}_{\mathcal{X}}} [\mathbf{x} \mathbf{x}^\top]$

Input: $\hat{\Sigma}, \text{Lev}_{\hat{\Sigma}, \mathcal{X}}$
repeat

$\mathbf{x}_1, \dots, \mathbf{x}_{d^2} \stackrel{\text{i.i.d.}}{\sim} \text{Lev}_{\hat{\Sigma}, \mathcal{X}}$

$\tilde{\mathbf{X}} \leftarrow \left[\frac{1}{\sqrt{d} l_{\hat{\Sigma}}(\mathbf{x}_i)} \cdot \mathbf{x}_i^\top \right]_{d^2 \times d}$

Sample $\text{Acc} \sim \text{Bernoulli} \left(\frac{\det(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})}{\det(\hat{\Sigma})} \right)$

until $\text{Acc} = \text{true}$

$S \leftarrow \text{VolSamp}(\tilde{\mathbf{X}})$

return \mathbf{X}_S

VolSamp($\mathbf{X} \in \mathbb{R}^{n \times d}$): [DW18]

$S \leftarrow \{1..n\}$

while $|S| > d$

$\forall i \in S \quad q_i \leftarrow \frac{\det(\mathbf{X}_{S \setminus i}^\top \mathbf{X}_{S \setminus i})}{(|S|-d) \det(\mathbf{X}_S^\top \mathbf{X}_S)}$

Sample $i \sim (q_i)_{i \in S}$

$S \leftarrow S \setminus \{i\}$

end

return S

DW18 Dereziński, Warmuth. *Reverse iterative volume sampling for linear regression*. JMLR, 2018

DWH18 Dereziński, Warmuth, Hsu. *Leveraged volume sampling for linear regression*. NeurIPS 2018

Der18 Dereziński. *Fast determinantal point processes via distortion-free intermediate sampling*. 2018