

Cramèr-Rao Bound for Estimation After Model Selection and its Application to Sparse Vector Estimation

Elad Meir^{1,2} and Tirza Cherlow-Routtenberg¹

[1] School of Electrical and Computer Engineering
Ben-Gurion University of the Negev, Israel.

[2] TechSee Augmented Vision Ltd.

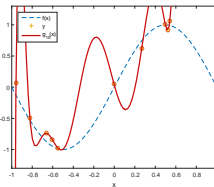
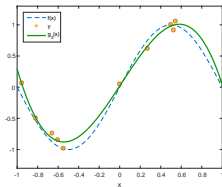
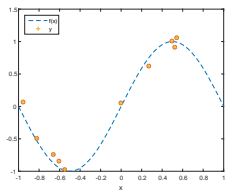


Overview

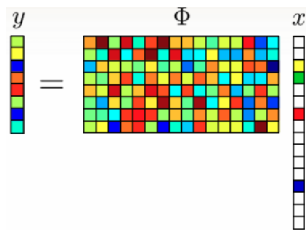
- 1 Motivation
- 2 Problem Formulation
- 3 Selective CRB
- 4 Examples
- 5 Our recent publications

Motivation

Coefficients estimation of polynomial regression



Sparse Signal processing



Estimation after model selection

The usual practice:

- select a model from a pool of candidate models, based on the data (e.g. AIC, MDL)
- estimate and analyze **selected** model, disregarding the selection process (maximum likelihood, least squares)

Applications:

multivariate data analysis, machine learning, graph analysis, ...

Non-Bayesian estimation after model selection

- **Goal:** include model selection process in the analysis and estimation
- **Previous works on the selection effect:**
 - effects on standard errors [Pöetscher 1991, Efron 2014]
 - confidence intervals [Benjamini and Yakutieli 2005, Kabaila and Leeb 2006]
 - Cramèr-Rao bound for signals under unknown model order (SMS) [Sando, Mitra and Stoica 2002].
 - Cramèr-Rao bound for estimation after *parameter selection* [Routtenberg and Tong 2016, Harel and Routtenberg 2019] - the model is known, selection of parameter of interest
- **Our approach:**
 - formulation of non-Bayesian post model selection estimation
 - coherence
 - (selective) unbiasedness
 - (selective) CRB

Outline

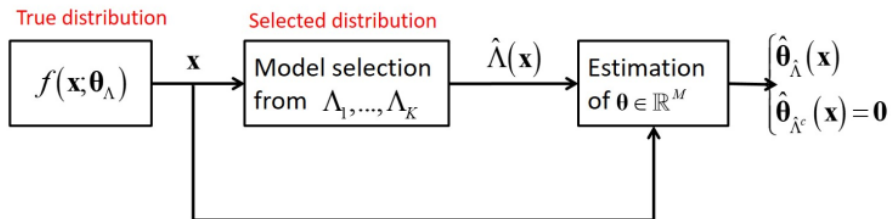
- 1 Motivation
- 2 Problem Formulation**
- 3 Selective CRB
- 4 Examples
- 5 Our recent publications

Problem formulation: estimation after model selection

- $\boldsymbol{\theta} = [\theta_1, \dots, \theta_M]^T \in \mathbb{R}^M$ - unknown deterministic parameter vector
- $\mathbf{x} \in \Omega_{\mathbf{x}}$ - random observation vector, which is **truly** distributed by
- $f(\mathbf{x}; \boldsymbol{\theta}_{\Lambda})$ - **true** probability density function
- Λ and Λ^c - the **true** support of $\boldsymbol{\theta}$ and its complement, according to the **true** model. s.t. $\boldsymbol{\theta}_{\Lambda} \in \mathbb{R}^{|\Lambda|}$, and $\boldsymbol{\theta}_{\Lambda^c} = \mathbf{0}$
- $\{\Lambda_k\}_k$ - is the set of all candidate supports, thus
- $\{f(\mathbf{x}; \boldsymbol{\theta}_{\Lambda_k})\}_k$ - set of pdfs, parameterized by $\boldsymbol{\theta}$ with the suitable support from $\{\Lambda_k\}_k$
- $\hat{\boldsymbol{\theta}} : \Omega_{\mathbf{x}} \rightarrow \mathbb{R}^M$ - estimator of $\boldsymbol{\theta}$ based on \mathbf{x}
- $\hat{\Lambda}$ - the **selection** rule
- $\pi_k(\boldsymbol{\theta}_{\Lambda}) \triangleq \Pr(\hat{\Lambda} = \Lambda_k; \boldsymbol{\theta}_{\Lambda})$ - the probability of **selecting** the k th model

Estimation after model selection:

The measurement vector, \mathbf{x} , is generated by the p.d.f., $f(\mathbf{x}; \boldsymbol{\theta}_\Lambda)$, then a model is selected by a pre-determined selection rule - $\hat{\Lambda}$. In the second stage, the unknown parameter, $\boldsymbol{\theta}$ is estimated based on the observation vector and the selected model.



Coherence estimation

Theorem

An estimator $\hat{\theta}$ is said to be a coherent estimator with respect to the selection rule $\hat{\Lambda}$, if $\hat{\theta}_{\hat{\Lambda}^c} = \mathbf{0}$

Structure:

- a model is selected according to a *predetermined* data-driven selection rule, $\hat{\Lambda}$
- the selected parameters, $\hat{\theta}_{\hat{\Lambda}}$, are estimated
- the deselected parameters are set to be zero, i.e. $\hat{\theta}_{\hat{\Lambda}^c} = \mathbf{0}$

Zero Padded (ZP) vectors and $\mathbf{D}_k(\Lambda)$ matrix

Definition

For an arbitrary vector, $\mathbf{a} \in \mathbb{R}^M$, and any candidate support set, Λ_k , $k = 1, \dots, K$, the vector $\mathbf{a}_{\Lambda_k}^{\text{ZP}}$, is a zero-padded, M -length vector, whose non-padded elements correspond to the elements of \mathbf{a}_{Λ_k} .

Definition

For any $k = 1, \dots, K$, the $\mathbf{D}_k(\Lambda)$ is a $M \times |\Lambda|$ matrix with the elements

$$[\mathbf{D}_k(\Lambda)]_{m,l} \triangleq \begin{cases} 1, & m \in \Lambda_k, m = [\Lambda]_l \\ 0, & \text{otherwise} \end{cases}, \quad (1)$$

$\forall l = 1, \dots, M$, where $[\Lambda]_l$ denotes the l th element of the true support set (not to be confused with the l th candidate support set, Λ_l).

Outline

- 1 Motivation
- 2 Problem Formulation
- 3 Selective CRB**
- 4 Examples
- 5 Our recent publications

cost function - motivation

- The Cramér-Rao bound (CRB) provides a lower bound on the mean squared error (MSE) of any mean-unbiased estimator and is used as a benchmark in non-Bayesian estimation
- The conventional CRB does not take into account the selection process, thus, it is inappropriate for estimation after model selection
- We wish to derive a CRB-type bound that takes into account the selection process

- The Cramér-Rao bound (CRB) provides a lower bound on the **mean squared error (MSE)** of any mean-unbiased estimator

The term mean square error (MSE) refers to the squared mean difference between the estimated and true values. Although this cost function is widely used in non-Bayesian estimation, it does not include the model selection procedure therefore neglects the selection process

Novel cost function

Selected-square-error (SSE) matrix, and the corresponding mean SSE (MSSE)

Proposition (selected-square-error matrix)

$$\mathbf{C}(\hat{\boldsymbol{\theta}}, \hat{\Lambda}, \boldsymbol{\theta}) \triangleq \left(\hat{\boldsymbol{\theta}}_{\hat{\Lambda}}^{ZP} - \boldsymbol{\theta}_{\hat{\Lambda}}^{ZP} \right) \left(\hat{\boldsymbol{\theta}}_{\hat{\Lambda}}^{ZP} - \boldsymbol{\theta}_{\hat{\Lambda}}^{ZP} \right)^T \quad (2)$$

Proposition (mean-selected-squared-error)

$$\mathbb{E}_{\boldsymbol{\theta}_{\Lambda}} \left[\mathbf{C}(\hat{\boldsymbol{\theta}}, \hat{\Lambda}, \boldsymbol{\theta}) \right] = \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_{\Lambda}) \mathbb{E}_{\boldsymbol{\theta}_{\Lambda}} \left[\left(\hat{\boldsymbol{\theta}}_{\Lambda_k}^{ZP} - \boldsymbol{\theta}_{\Lambda_k}^{ZP} \right) \left(\hat{\boldsymbol{\theta}}_{\Lambda_k}^{ZP} - \boldsymbol{\theta}_{\Lambda_k}^{ZP} \right)^T \mid \hat{\Lambda} = \Lambda_k \right] \quad (3)$$

A scalar version is available in the article

- The Cramér-Rao bound (CRB) provides a lower bound on the mean squared error (MSE) of any **mean-unbiased** estimator

Mean-unbiasedness (*unbiasedness*) considers the estimated, and expected values, but does not include the model selection procedure.

Moreover, most post selection estimator tend to bias, making the CRB unsuitable for our use case.

Proposition (Selective Unbiasdness)

An estimator, $\hat{\boldsymbol{\theta}}$, is an unbiased estimator for the problem of estimating the true parameter vector, $\boldsymbol{\theta}_\Lambda$, in the Lehmann sense w.r.t. the SSE matrix defined earlier, and the selection rule, $\hat{\Lambda}$, iff

$$\mathbf{b}_k(\boldsymbol{\theta}, \Lambda) \triangleq \mathbb{E}_{\boldsymbol{\theta}_\Lambda} \left[\hat{\boldsymbol{\theta}}_{\Lambda_k}^{ZP} - \boldsymbol{\theta}_{\Lambda_k}^{ZP} \mid \hat{\Lambda} = \Lambda_k \right] = \mathbf{0}, \forall \boldsymbol{\theta}_{\Lambda_k}^{ZP} \in \mathbb{R}^{|\Lambda|}, \quad (4)$$

for all $k = 1, \dots, K$, such that $\pi_k(\boldsymbol{\theta}_\Lambda) \neq 0$.

proof appears in the article, Appendix A

In practice, coherent post model selection estimators tend to be biased.
Thus (4) is not zero

Proposition (Selective bias)

$$\mathbf{b}_k(\boldsymbol{\theta}, \Lambda) \triangleq \mathbb{E}_{\boldsymbol{\theta}_\Lambda} \left[\hat{\boldsymbol{\theta}}_{\Lambda_k}^{ZP} - \boldsymbol{\theta}_{\Lambda_k}^{ZP} \mid \hat{\Lambda} = \Lambda_k \right] \neq \mathbf{0}, \forall k = 1, \dots, K, \quad (5)$$

while its derivative w.r.t. the true parameter is defined as

$$\mathbf{G}_k(\boldsymbol{\theta}, \Lambda) \triangleq \nabla_{\boldsymbol{\theta}_\Lambda} \mathbf{b}_k(\boldsymbol{\theta}_\Lambda), \quad \forall k = 1, \dots, K. \quad (6)$$

The selection bias affects the relation between the MSE and the MSSE of the estimator, as we can see in

Proposition (MSE-MSSE coherent relations)

For any *coherent* estimator, the MSE satisfies

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}, \Lambda) &= \mathbb{E}_{\boldsymbol{\theta}_\Lambda} \left[\mathbf{C}(\hat{\boldsymbol{\theta}}, \hat{\Lambda}, \boldsymbol{\theta}) \right] + \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_\Lambda) \boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} (\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}})^T \\ &\quad - \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_\Lambda) \left(\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} \mathbf{b}_k^T(\boldsymbol{\theta}, \Lambda) + \mathbf{b}_k(\boldsymbol{\theta}, \Lambda) (\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}})^T \right). \end{aligned} \quad (7)$$

proof appears in the article, Appendix B

Proposition (Selective CRB on the MSSE)

Under regularity conditions, the MSSE of any coherent and selective biased estimator, with the selective bias from (4), is bounded by

$$\mathbb{E}_{\boldsymbol{\theta}_\Lambda}[\mathbf{C}(\hat{\boldsymbol{\theta}}, \hat{\Lambda}, \boldsymbol{\theta}_\Lambda)] \succeq \mathbf{B}_{sCRB}(\boldsymbol{\theta}_\Lambda), \quad (8)$$

where the biased selective CRB is defined as

$$\mathbf{B}_{sCRB}(\boldsymbol{\theta}_\Lambda) \triangleq \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_\Lambda) [(\mathbf{D}_{\Lambda_k}^{ZP} + \mathbf{G}_{\Lambda_k}^{ZP}) \mathbf{J}_k^{-1} (\mathbf{D}_{\Lambda_k}^{ZP} + \mathbf{G}_{\Lambda_k}^{ZP})^T] \quad (9)$$

selective CRB

The selective CRB matrix bound

Theorem (sCRB on the MSE)

For any *coherent* estimator, the MSE satisfies

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}_\Lambda) \succeq & \mathbf{B}_{\text{sCRB}}(\boldsymbol{\theta}_\Lambda) + \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_\Lambda) \boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} \left(\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} \right)^T \\ & - \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_\Lambda) \left(\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} \mathbf{b}_k^T(\boldsymbol{\theta}, \Lambda) + \mathbf{b}_k(\boldsymbol{\theta}, \Lambda) (\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}})^T \right). \end{aligned}$$

Proof (Appendix C), and the marginal version appear in the article (Section C)

selective unbiased CRB

The selective CRB matrix bound

Theorem (sCRB on the MSE)

For any *coherent* and *selective unbiased* estimator, the MSE satisfies

$$\text{MSE}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}_\Lambda) \succeq \mathbf{B}_{\text{sCRB}}(\boldsymbol{\theta}_\Lambda) + \sum_{k=1}^K \pi_k(\boldsymbol{\theta}_\Lambda) \boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} \left(\boldsymbol{\theta}_{\Lambda_k^c}^{\text{ZP}} \right)^T.$$

Outline

- 1 Motivation
- 2 Problem Formulation
- 3 Selective CRB
- 4 Examples**
- 5 Our recent publications

General Linear Model

The General Linear Model (GLM) is given by

- its observations $\forall k = 1, \dots, K$,

$$\mathbf{x} = \mathbf{H}_k \boldsymbol{\theta}_{\Lambda_k} + \mathbf{w}, \quad (10)$$

- the matrices $\mathbf{H}_k \in \mathbb{R}^{N \times |\Lambda_k|}$, $k = 1, \dots, K$, are assumed to be known full column rank matrices
- $\boldsymbol{\theta} \in \mathbb{R}^M$ is a deterministic unknown vector
- \mathbf{w} is a zero-mean i.i.d. Gaussian random vector, with known variance
- the coherent Maximum Selected Likelihood (MSL) estimator, is given for $k = 1, \dots, K$, by

$$\hat{\boldsymbol{\theta}}_{\Lambda_k}^{\text{ML}|k} = (\mathbf{H}_k^T \mathbf{H}_k)^{-1} \mathbf{H}_k^T \mathbf{x}, \quad (11)$$

- coherency dictates $\hat{\boldsymbol{\theta}}_{\Lambda_k^c} = \mathbf{0}$

Vs. SNR

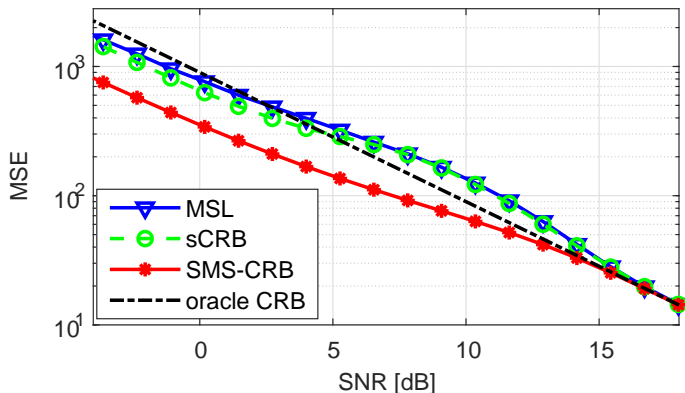


Figure: General linear model with AIC selection rule: MSE, sCRB, SMS-CRB, oracle CRB, vs. $\text{SNR} \triangleq 10 \log_{10} \frac{\|\mathbf{H}\boldsymbol{\theta}\|^2}{N\sigma^2}$, with varying σ , $N = 1500$ samples, $\boldsymbol{\theta} = [4, -3]^T$, $\mathbf{h}_1 = [1, \dots, 1]^T$, and the values of \mathbf{h}_2 are randomly drawn from uniform distribution in interval $[0, 10]$.

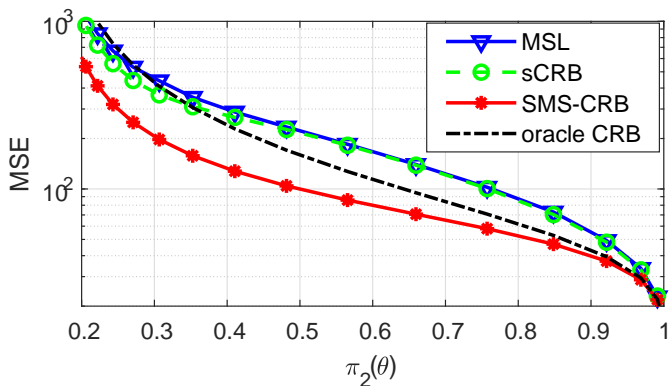
Vs. π_k 

Figure: General linear model with AIC selection rule: MSE, sCRB, SMS-CRB, oracle CRB, vs. $\text{SNR} \triangleq 10 \log_{10} \frac{\|\mathbf{H}\boldsymbol{\theta}\|^2}{N\sigma^2}$, with varying σ , $N = 1500$ samples, $\boldsymbol{\theta} = [4, -3]^T$, $\mathbf{h}_1 = [1, \dots, 1]^T$, and the values of \mathbf{h}_2 are randomly drawn from uniform distribution in interval $[0, 10]$.

Vs. GIC functions

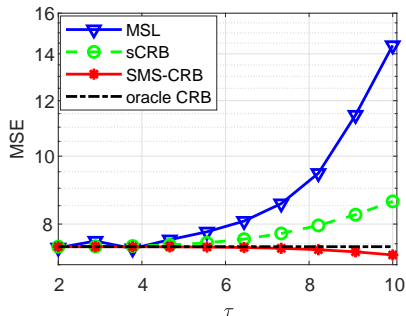
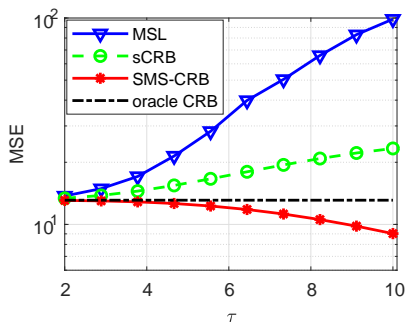


Figure: General linear model with GIC selection rule: The MSE of the MSL estimator compared to selective CRB and the SMS-CRB versus different values of the parameter $\tau(N, |\Lambda_k|)$, $N = 150$, with SNR = -3.5dB (left) and 0dB (right).

OST

For sparse vector estimation, we inspect the MSE of

$$\hat{\boldsymbol{\theta}}_{\hat{\Lambda}}^{\text{ML-OST}} = \left(\mathbf{A}_{\hat{\Lambda}}^T \mathbf{A}_{\hat{\Lambda}} \right)^{-1} \mathbf{A}_{\hat{\Lambda}}^T \mathbf{x}, \quad (12)$$

where the selection function is OST

$$m \in \hat{\Lambda} \text{ if } |\mathbf{a}_m^T \mathbf{x}| > c > 0, \quad \forall m = 1, \dots, M$$

and, of course, coherency dictates $\hat{\boldsymbol{\theta}}_{\hat{\Lambda}^c}^{\text{ML-OST}} = \mathbf{0}$.

1st example

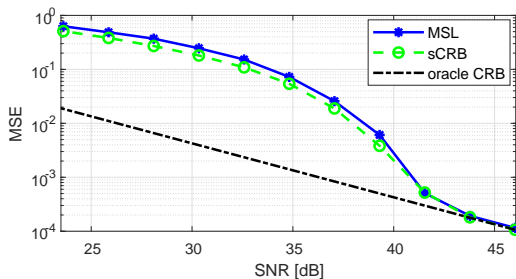


Figure: $\theta_{\Lambda} = 1$. The MSE of the MSL estimator, sCRB, and oracle CRB under OST rule, with $c = 0.95$, versus SNR. Random 7×14 dictionary, \mathbf{A} , $\mathbf{a}_m^T \mathbf{a}_m = 1$, mutual coherence is $\mu = 0.5673$. $|\Lambda| = 3$

2nd example

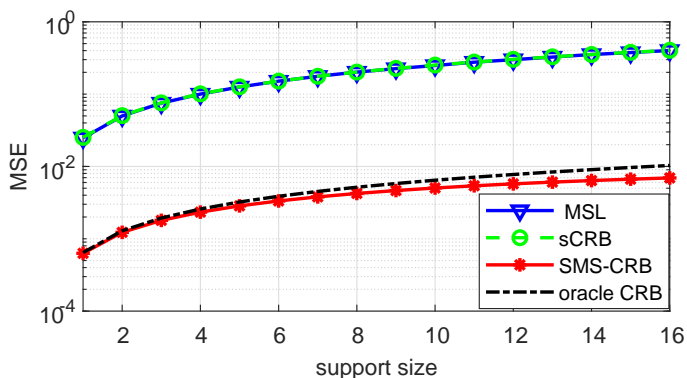


Figure: $\theta_{\Lambda} = 1$. The MSE of the MSL estimator, sCRB, and oracle CRB under OST rule, with $c = 0.95$. Hadamard 16×16 dictionary, \mathbf{A} , $\mathbf{a}_m^T \mathbf{a}_m = 1$, $\sigma = 0.1594$, $1 \leq |\Lambda| \leq 16$

Biased sCRB with identity dictionary, $\mathbf{A} = \mathbf{I}$, with $L = M$.

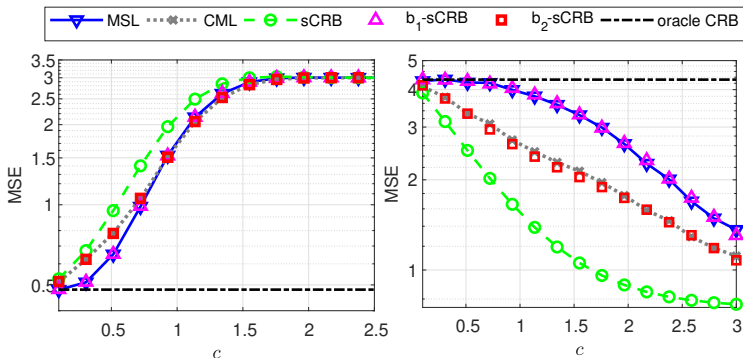


Figure: Sparse vector estimation: The MSE of the MSL and of the post-selection conditional ML (CML) estimators, where the likelihood is selected by the OST rule, the selective CRB, the biased selective CRB with the MSL bias, b_1 -sCRB, the biased selective CRB with the CML bias, b_2 -sCRB, and the oracle CRB, versus the threshold, c , for 1) $\theta_m = 1, \sigma = 0.4$ (left); 2) $\theta_m = 0.5, \sigma = 1.2$ (right).

Conclusion

- formulation of **coherent** estimation after model selection
- a measure of performance - **selected-square-error** is considered
- a Lehmann-sense **selective unbiasedness** definition is introduced
- an appropriate **selective Cramér Rao Bound** is derived, with its biased version and for sparse vector estimation
- simulations show that the selective CRB is valid, and tighter than the SMS-CRB, and predicts the threshold phenomenon



Outline

- 1 Motivation
- 2 Problem Formulation
- 3 Selective CRB
- 4 Examples
- 5 Our recent publications**

Our recent publications:

- N. Harel and T. Routtenberg, "Post-Parameter-Selection Maximum-Likelihood Estimation", IEEE Workshop on Statistical Signal Processing (SSP 2021).
- E. Meir and T. Routtenberg, "Cramér-Rao Bound for Estimation After Model Selection with Applications to Sparse Vector Estimation", IEEE Transactions on Signal Processing", vol. 69, pp. 2284-2301, 2021.
- N. Harel and T. Routtenberg, "Bayesian Estimation After Model Selection", IEEE Signal Processing Letters, vol. 28, pp. 175-179, 2021.
- T. Weiss, T. Routtenberg, and H. Messer, "Total Performance Evaluation of Intensity Estimation after Detection, Signal Processing, vol. 183, pp. 1-8, 2021.
- N. Harel and T. Routtenberg, "Low-Complexity Methods for Estimation After Parameter Selection, IEEE Transactions on Signal Processing, vol. 68, pp. 1152-1167, 2020.
- E. Meir and T. Routtenberg, "Selective Cramr-Rao Bound for Estimation After Model Selection", IEEE Workshop on Statistical Signal Processing (SSP 2018).
- S. Cohen, T. Routtenberg, and L. Tong, "Non-Bayesian Parameter Estimation of the Probability of the Missing Mass", arXiv:2101.04329.

