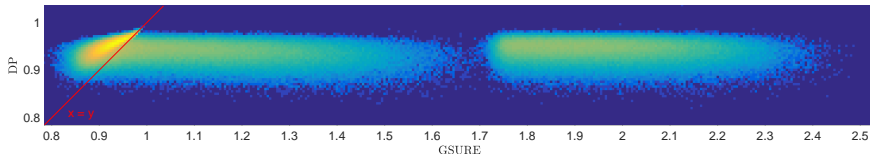


## The Ill-Posedness Always Rings Twice

### Risk Estimators for Choosing Regularization Parameters in Inverse Problems



**Felix Lucka**, University College London, [f.lucka@ucl.ac.uk](mailto:f.lucka@ucl.ac.uk)

**joint with:** Katharina Proksch, Christoph Brune, Nicolai Bissantz, Martin Burger, Holger Dette & Frank Wübbeling

Discrete inverse problem:

$$y = Ax^* + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_m)$$

Variational regularization:

$$\hat{x}_\alpha(y) = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \|Ax - y\|_2^2 + \alpha R(x),$$

$R$  convex such that the minimizer is unique for  $\alpha > 0$ .

Discrete inverse problem:

$$y = Ax^* + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_m)$$

Variational regularization:

$$\hat{x}_\alpha(y) = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \|Ax - y\|_2^2 + \alpha R(x),$$

$R$  convex such that the minimizer is unique for  $\alpha > 0$ .

Every talk: "How did you choose  $\alpha$ ?"

A problem as old as inverse problems / robust statistical inference.

- ▶ Ideal parameter choice:

$$\alpha^* := \operatorname{argmin}_{\alpha \geq 0} \|\hat{x}_\alpha(y) - x^*\|_2^2$$

! obviously not available (**oracle solution**)

- ▶ Many different approaches proposed. Focus here: Strategies that need accurate estimate of noise variance  $\sigma^2$ .
- ▶ Classical example: **discrepancy principle**:

$$\text{find } \alpha \text{ s.t. } \|A\hat{x}_\alpha(y) - y\|_2^2 = m\sigma^2.$$

- ✓ robust and easy-to-implement for many applications
- ! typically over-estimates  $\alpha^*$ .

We want to minimize the quadratic risk function

$$R_{\text{SURE}}(\alpha) := \mathbb{E} [\|Ax^* - A\hat{x}_\alpha(y)\|_2^2],$$

but as  $R_{\text{SURE}}$  depends on  $x^*$ , we replace it by an **unbiased estimate**:

$$\text{SURE}(\alpha, y) := \|y - A\hat{x}_\alpha(y)\|_2^2 - m\sigma^2 + 2\sigma^2 \text{df}_\alpha(y), \quad \text{df}_\alpha(y) = \text{tr}(\nabla_y \cdot A\hat{x}_\alpha(y)),$$

where unbiased means:  $\mathbb{E}[\text{SURE}(\alpha, y)] = R_{\text{SURE}}(\alpha)$

Risk in the domain, not in the image of the operator  $A$ :

$$R_{\text{GSURE}}(\alpha) := \mathbb{E} [\|\Pi(x^* - \hat{x}_\alpha(y))\|_2^2], \quad \Pi := A^+A$$

$$\text{GSURE}(\alpha, y) := \|x_{\text{ML}}(y) - \hat{x}_\alpha(y)\|_2^2 - \sigma^2 \text{tr}((AA^*)^+) + 2\sigma^2 \text{gdf}_\alpha(y)$$

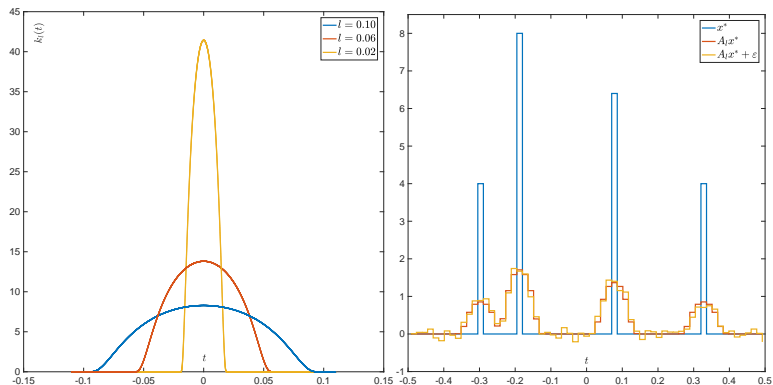
$$\text{gdf}_\alpha(y) := \text{tr}((AA^*)^+ \nabla_y A \hat{x}_\alpha(y)), \quad x_{\text{ML}} = A^+y = A^*(AA^*)^+y,$$

**GSURE** seems more appropriate for ill-posed problems, since properties in data space do not tell much about the reconstruction quality!

Recently, a lot of work on risk estimators in imaging and inverse problems:  
*Blu, Chesneau, Deledalle, Dossal, Elad, Eldar, Fadili, Giryes, Kachour, Kocher, Luisier, Morel, Peyré, Ramani, Unser, Vaiter, Van De Ville, Wang*

Our interest is a statistical perspective:

- ▶ All parameter choice rules depend on data  $y$  and hence on random  $\varepsilon$ .
- ▶ Therefore,  $\hat{\alpha}_{DP}$ ,  $\hat{\alpha}_{SURE}$  and  $\hat{\alpha}_{GSURE}$  are **random variables**.
- ▶ Characteristics of their **probability distributions**?
- ▶ **Distributions or error measures**  $dist(x^*, x_{\hat{\alpha}})$ ?



$$y_\infty(s) = A_{\infty, l} x_\infty^*, \quad x_\infty^*(t) := \sum_{i=1}^4 a_i \delta(b_i)$$

1D periodic convolution, kernel width  $l$ , mass-preserving discretization into ONB of piecewise constant functions

$$y_m = A_{m, l} x_m^* + \varepsilon_m, \quad \varepsilon_m \sim \mathcal{N}(0, \sigma^2 l_m)$$

Quadratic regularization leading to explicit, linear estimator:

$$\hat{x}_\alpha(y) = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \|Ax - y\|_2^2 + \frac{\alpha}{2} \|x\|_2^2 = (A^*A + \alpha I)^{-1} A^*y$$

Switch to singular system to analyse:

$$A = U\Sigma V^*, \quad 1 = \gamma_1 \geq \dots \geq \gamma_m > 0$$

$$y_i = \langle u_i, y \rangle, \quad x_i^* = \langle v_i, x^* \rangle, \quad \tilde{\epsilon}_i = \langle u_i, \epsilon \rangle$$

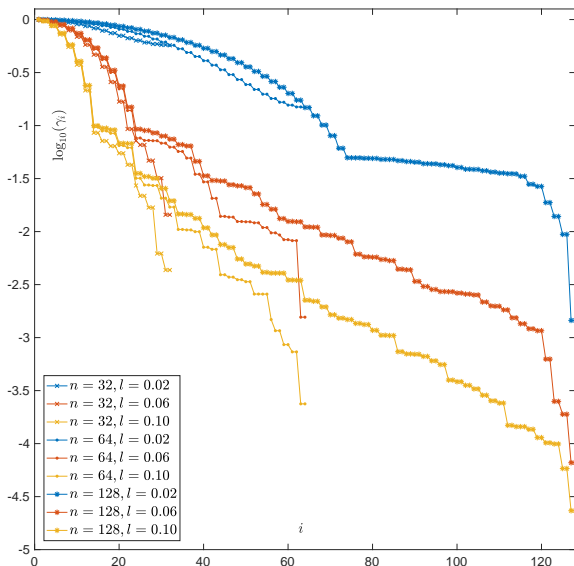
$$y = Ax + \epsilon \Leftrightarrow y_i = \gamma_i x_i^* + \tilde{\epsilon}_i, \quad \tilde{\epsilon}_i \sim \mathcal{N}(0, \sigma^2)$$

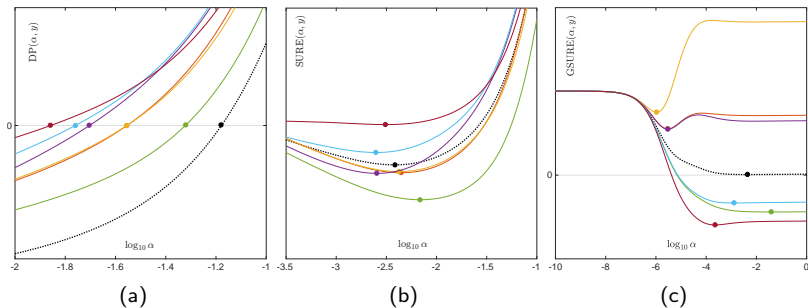


$$\begin{aligned} \text{DP}(\alpha, y) &:= \|A\hat{x}_\alpha(y) - y\|_2^2 - m\sigma^2 \\ &= \sum_{i=1}^m \frac{\alpha^2}{(\gamma_i^2 + \alpha)^2} y_i^2 - m\sigma^2 \end{aligned}$$

$$\begin{aligned} \text{SURE}(\alpha, y) &= \|y - A\hat{x}_\alpha(y)\|_2^2 - m\sigma^2 + 2\sigma^2 \text{df}_\alpha(y) \\ &= \sum_{i=1}^m \frac{\alpha^2}{(\gamma_i^2 + \alpha)^2} y_i^2 - m\sigma^2 + 2\sigma^2 \sum_{i=1}^m \frac{\gamma_i^2}{\gamma_i^2 + \alpha} \end{aligned}$$

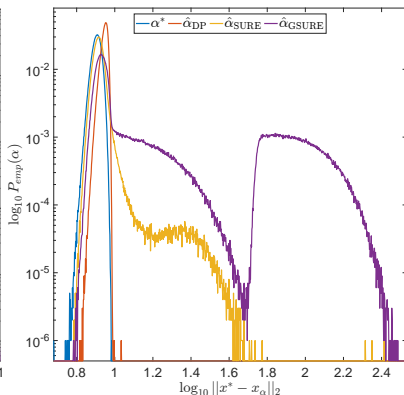
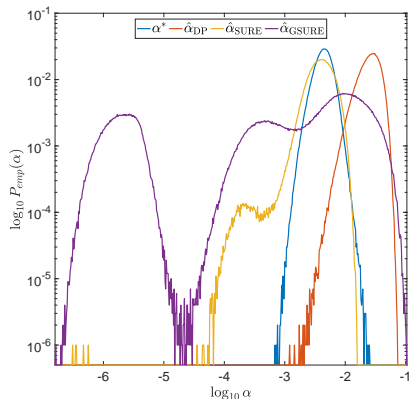
$$\begin{aligned} \text{GSURE}(\alpha, y) &= \|x_{\text{ML}}(y) - \hat{x}_\alpha(y)\|_2^2 - \sigma^2 \text{tr}((AA^*)^+) + 2\sigma^2 \text{gdf}_\alpha(y) \\ &= \sum_{i=1}^r \left( \frac{1}{\gamma_i} - \frac{\gamma_i}{\gamma_i^2 + \alpha} \right)^2 y_i^2 - \sigma^2 \sum_{i=1}^r \frac{1}{\gamma_i^2} + 2\sigma^2 \sum_{i=1}^r \frac{1}{\gamma_i^2 + \alpha} \end{aligned}$$

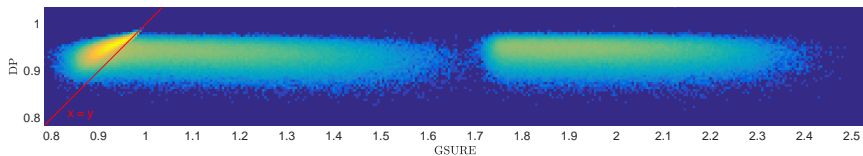
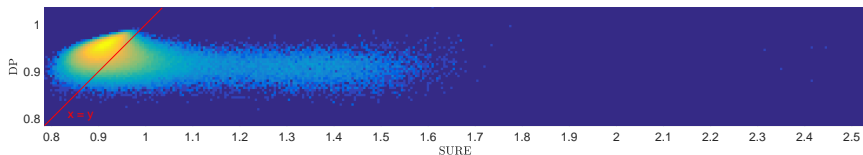
Singular values  $\gamma_i$  of  $A_l$  for different choices of  $m$  and  $l$ .



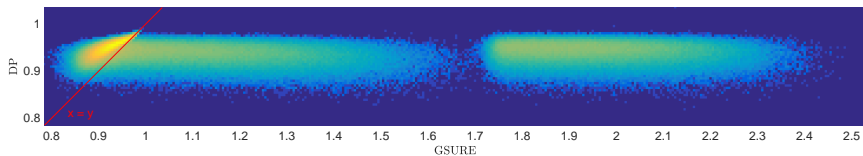
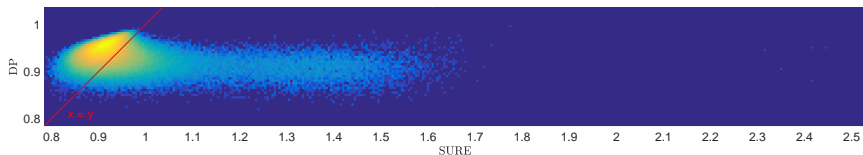
- (a)  $R_{DP}(\alpha) = \|A\hat{x}_\alpha(Ax^*) - Ax^*\|_2^2 - m\sigma^2$  vs. 6 realizations of  $DP(\alpha, y) = \|A\hat{x}_\alpha(y) - y\|_2^2 - m\sigma^2$
- (b)  $R_{SURE}(\alpha) = \mathbb{E} [\|Ax^* - A\hat{x}_\alpha(y)\|_2^2]$  vs. 6 realizations of  $SURE(\alpha, y) = \|y - A\hat{x}_\alpha(y)\|_2^2 - m\sigma^2 + 2\sigma^2 df_\alpha(y)$ .
- (c)  $R_{GSURE}(\alpha) = \mathbb{E} [\|\Pi(x^* - \hat{x}_\alpha(y))\|_2^2]$  vs. 6 realizations of  $GSURE(\alpha, y) = \|x_{ML}(y) - \hat{x}_\alpha(y)\|_2^2 - \sigma^2 \text{tr}((AA^*)^+) + 2\sigma^2 gdf_\alpha(y)$

- ▶ fine logarithmical  $\alpha$ -grid:  $\log_{10}(\alpha_i)$  from  $-40$  to  $40$ , step size  $0.01$ .
- ▶  $N_\varepsilon = 10^6$  samples of  $\varepsilon$ .
- ▶  $m = n = 64$ ,  $l = 0.06$ ,  $\sigma = 0.1$



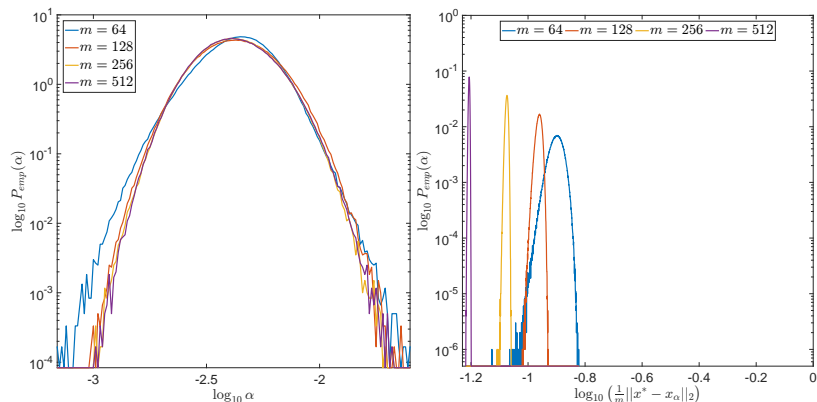


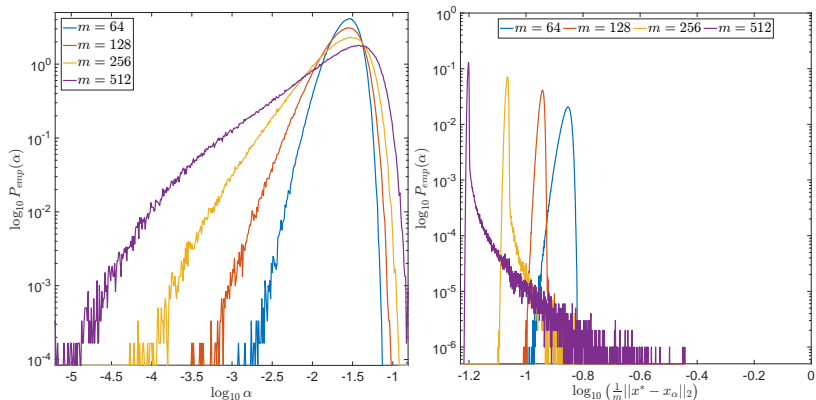
Joint empirical log-probabilities of  $\log_{10} \|x^* - x_{\hat{\alpha}}\|_2$



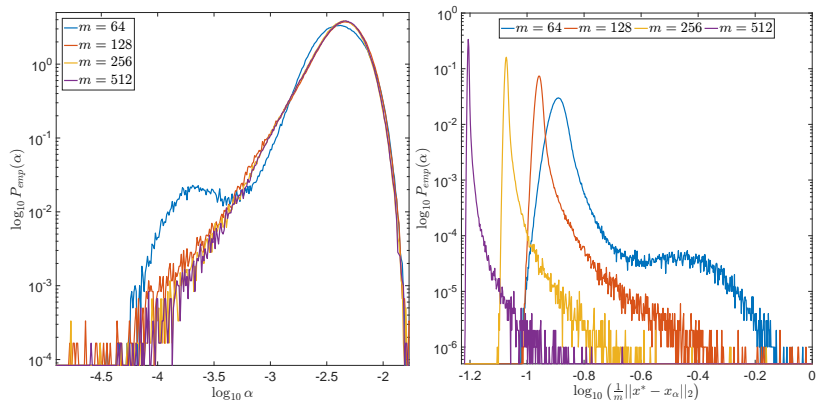
Joint empirical log-probabilities of  $\log_{10} \|x^* - x_{\hat{\alpha}}\|_2$

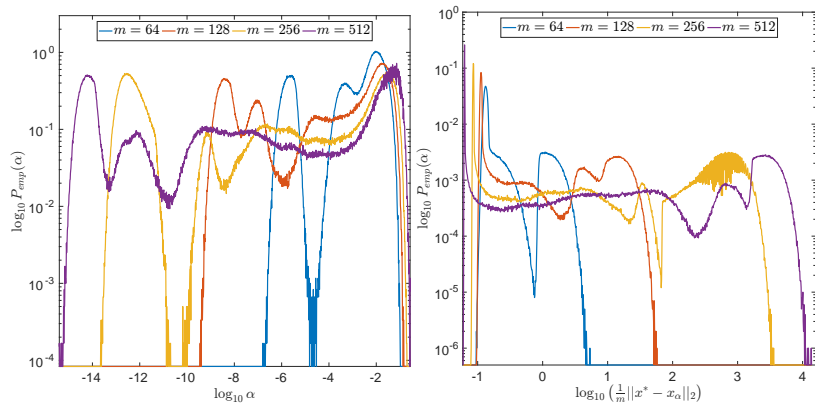
What's wrong? Let's do some more numerical studies...

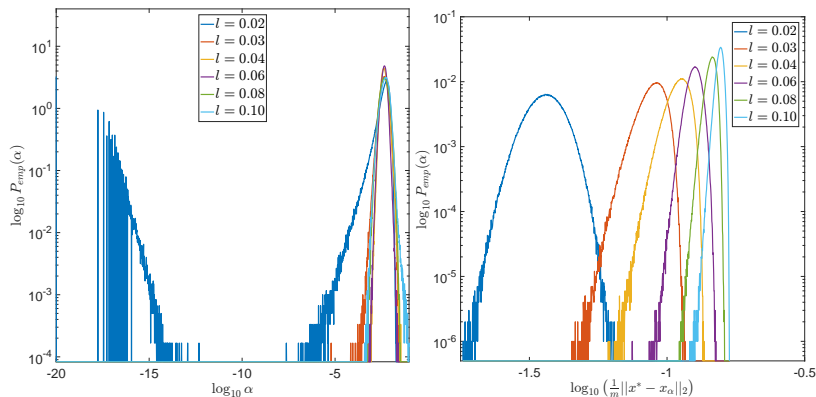


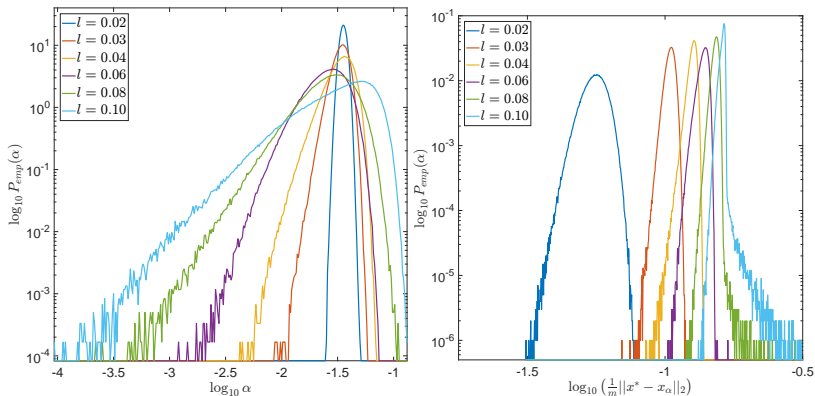


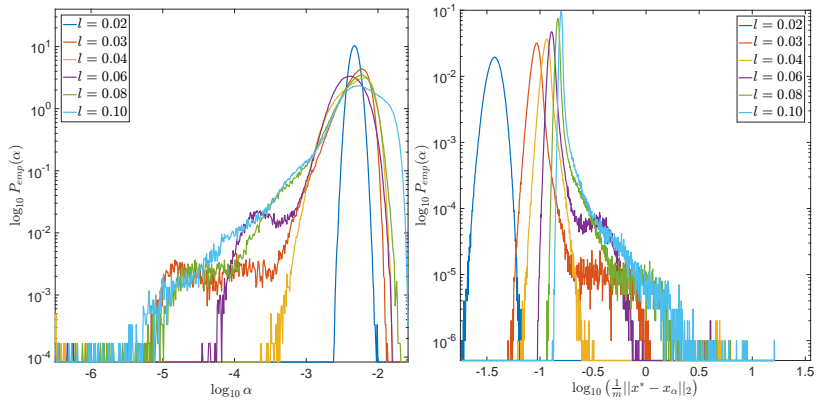


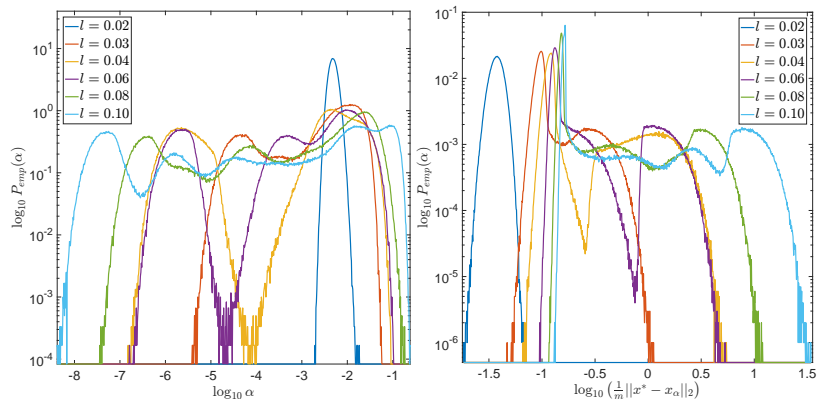












Assume  $A \in R^{m \times m}$ ,  $1 = \gamma_1 \geq \dots \geq \gamma_m > 0$  and  $\|x^*\|_2^2 = O(m)$ .

One can prove that for  $m \rightarrow \infty$ :

$$\sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{SURE}(\alpha, y) - R_{\text{SURE}}(\alpha, y)) \right| = O_{\mathbb{P}} \left( \frac{1}{\sqrt{m}} \right)$$

$$\sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{DP}(\alpha, y) - \mathbb{E}(\text{DP}(\alpha, y))) \right| = O_{\mathbb{P}} \left( \frac{1}{\sqrt{m}} \right)$$

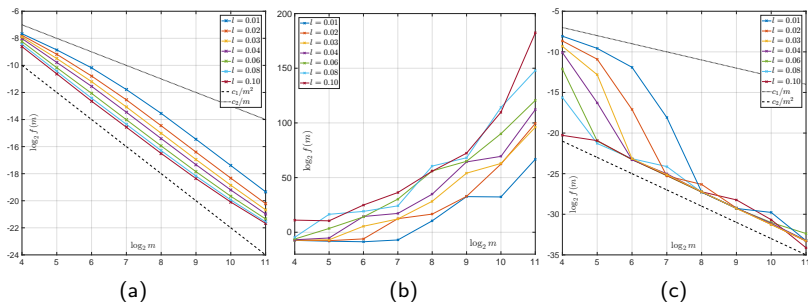
$$\sup_{\alpha \in [0, \infty)} \left| \frac{1}{m \text{cond}(A_m)^2} (\text{GSURE}(\alpha, y) - R_{\text{GSURE}}(\alpha)) \right| = O_{\mathbb{P}} \left( \frac{1}{\sqrt{m}} \right)$$

$$\mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{SURE}(\alpha, y) - R_{\text{SURE}}(\alpha, y)) \right| \right)^2 = O \left( \frac{1}{m} \right)$$

$$\mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{DP}(\alpha, y) - \mathbb{E}(\text{DP}(\alpha, y))) \right| \right)^2 = O \left( \frac{1}{m} \right)$$

$$\mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m \text{cond}(A_m)^2} (\text{GSURE}(\alpha, y) - R_{\text{GSURE}}(\alpha)) \right| \right)^2 = O \left( \frac{1}{m} \right)$$

**Proof:** Kolmogorov's maximal inequality & Doob's martingale inequality.



$$(a) \mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{SURE}(\alpha, y) - R_{\text{SURE}}(\alpha, y)) \right| \right)^2$$

$$(b) \mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{GSURE}(\alpha, y) - R_{\text{GSURE}}(\alpha)) \right| \right)^2$$

$$(c) \mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m \text{cond}(A_m)^2} (\text{GSURE}(\alpha, y) - R_{\text{GSURE}}(\alpha)) \right| \right)^2$$






Sparsity-inducing regularization (**LASSO**):

$$\hat{x}_\alpha(y) = \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - y\|_2^2 + \alpha \|x\|_1 \quad (1)$$

Let  $I$  be the support of  $\hat{x}_\alpha(y)$ ,  $|I| = k$ ,  $P_I \in \mathbb{R}^{k \times n}$  projector onto  $I$ ,  $A_I$  restriction of  $A$  to  $I$ . For our setting, we have that

$$\operatorname{df}_\alpha = \|\hat{x}_\alpha(y)\|_0 = k, \quad \operatorname{gdf}_\alpha = \operatorname{tr}(\Pi P_I (A_I^* A_I)^{-1} P_I^*)$$

-  **Deledalle, Vaïter, Peyré, Fadili, Dossal, 2012.** *Unbiased risk estimation for sparse analysis regularization*, [IEEE ICIP](#).
-  **Vaïter, Deledalle, Peyré, Fadili, Dossal, 2014.** *The Degrees of Freedom of Partly Smooth Regularizers*, [arXiv:1404.5557](#).
-  **Vaïter, Deledalle, Peyré, Dossal, Fadili, 2013.** *Local behavior of sparse analysis regularization: Applications to risk estimation*, [Applied and Computational Harmonic Analysis 35\(3\)](#).

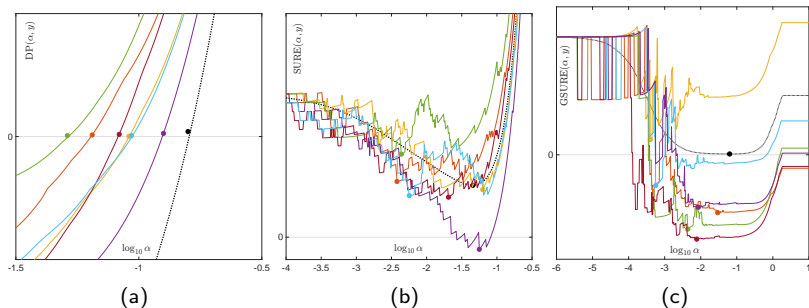
Sparsity-inducing regularization (LASSO):

$$\hat{x}_\alpha(y) = \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - y\|_2^2 + \alpha \|x\|_1 \quad (1)$$

Let  $I$  be the support of  $\hat{x}_\alpha(y)$ ,  $|I| = k$ ,  $P_I \in \mathbb{R}^{k \times n}$  projector onto  $I$ ,  $A_I$  restriction of  $A$  to  $I$ . For our setting, we have that

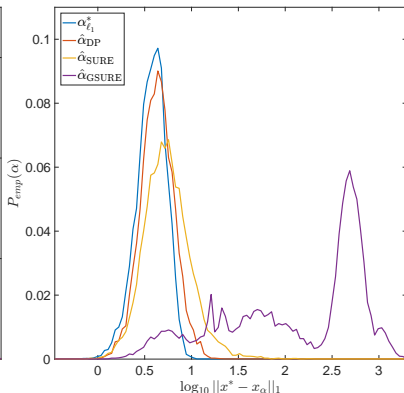
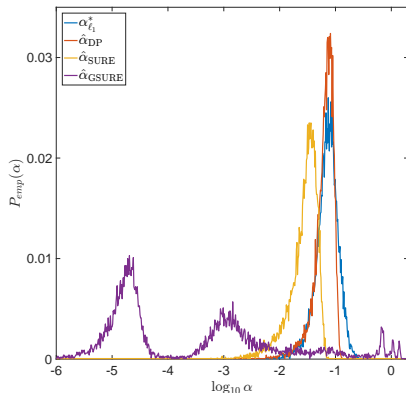
$$\operatorname{df}_\alpha = \|\hat{x}_\alpha(y)\|_0 = k, \quad \operatorname{gdf}_\alpha = \operatorname{tr}(\Pi P_I (A_I^* A_I)^{-1} P_I^*)$$

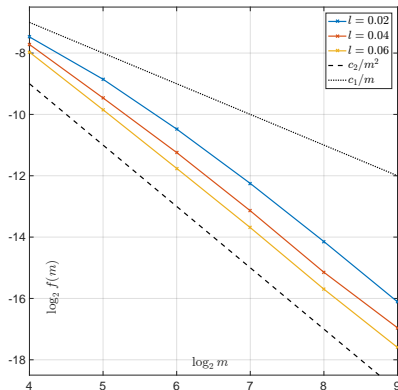
- ! No theory, only numerical studies.
- ! Fast but accurate and consistent computation of  $\hat{x}_\alpha(y)$  for  $\alpha$ 's ranging from  $10^{-10}$  to  $10^{10}$ .
- ✓ all-at-once implementation of ADMM solving (1) for all  $\alpha$  simultaneously with  $\operatorname{tol} = 10^{-14}$  and  $10^4$  max iter.



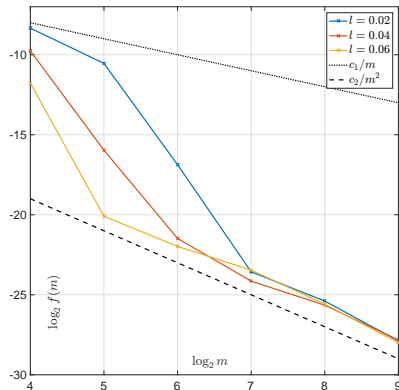
- (a)  $R_{DP}(\alpha) = \|\hat{A}\hat{x}_\alpha(Ax^*) - Ax^*\|_2^2 - m\sigma^2$  vs. 6 realizations of  $DP(\alpha, y) = \|A\hat{x}_\alpha(y) - y\|_2^2 - m\sigma^2$
- (b)  $R_{SURE}(\alpha) = \mathbb{E} [\|Ax^* - A\hat{x}_\alpha(y)\|_2^2]$  vs. 6 realizations of  $SURE(\alpha, y) = \|y - A\hat{x}_\alpha(y)\|_2^2 - m\sigma^2 + 2\sigma^2 df_\alpha(y)$ .
- (c)  $R_{GSURE}(\alpha) = \mathbb{E} [\|\Pi(x^* - \hat{x}_\alpha(y))\|_2^2]$  vs. 6 realizations of  $GSURE(\alpha, y) = \|x_{ML}(y) - \hat{x}_\alpha(y)\|_2^2 - \sigma^2 \text{tr}((AA^*)^+) + 2\sigma^2 \text{gdf}_\alpha(y)$

- ▶ fine logarithmical  $\alpha$ -grid:  $\log_{10}(\alpha_i)$  from  $-10$  to  $10$ , step size  $0.01$ .
- ▶  $N_\varepsilon = 10^4$  samples of  $\varepsilon$ .
- ▶  $m = n = 64$ ,  $l = 0.06$ ,  $\sigma = 0.1$





(a)

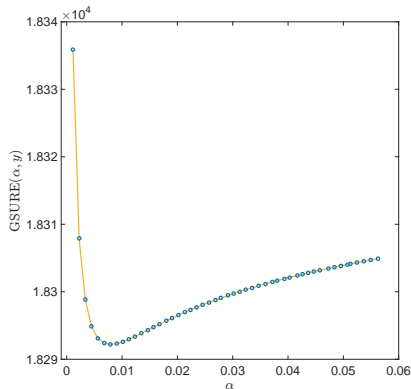


(b)

$$(a) \mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m} (\text{SURE}(\alpha, y) - R_{\text{SURE}}(\alpha, y)) \right| \right)^2$$

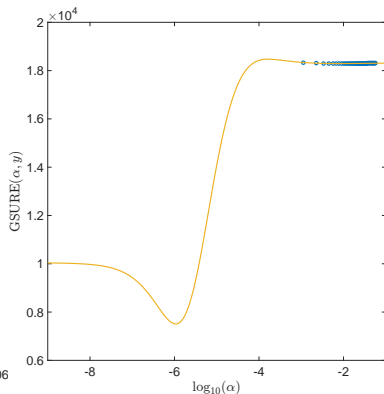
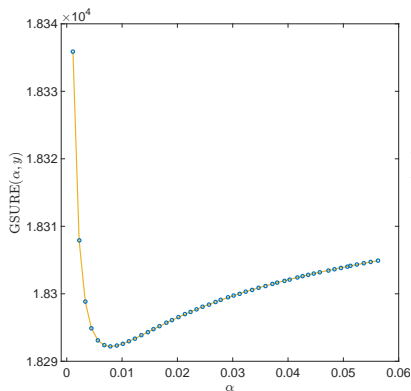
$$(b) \mathbb{E} \left( \sup_{\alpha \in [0, \infty)} \left| \frac{1}{m \text{cond}(A_m)^2} (\text{GSURE}(\alpha, y) - R_{\text{GSURE}}(\alpha)) \right| \right)^2$$

GSURE computed on a linear grid around "a reasonable value"...



(quadratic regularization)

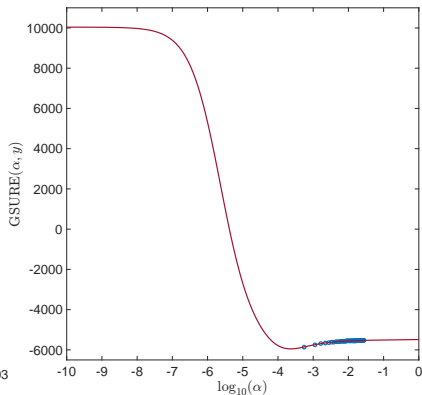
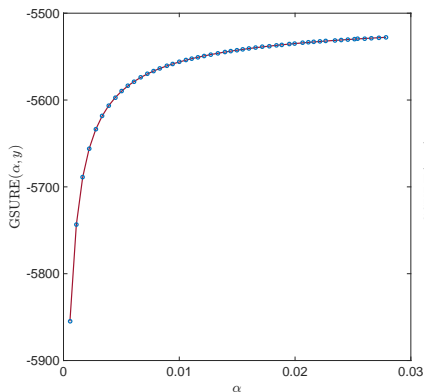
GSURE computed on a linear grid around "a reasonable value"...



...and on a fine logarithmic grid.

(quadratic regularization)

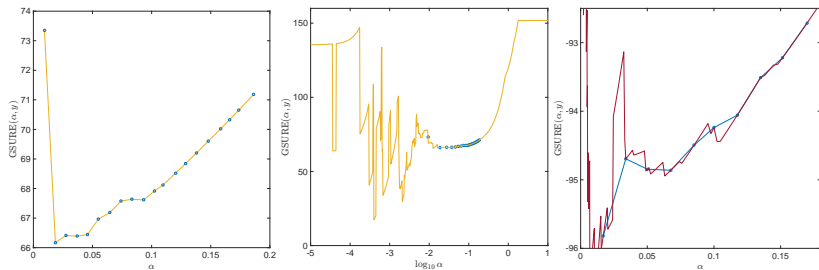
GSURE computed on a linear grid around "a reasonable value"...



...and on a fine logarithmic grid.

(quadratic regularization)

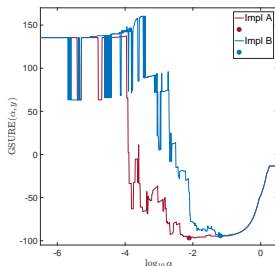




(LASSO regularization)

In addition to fine logarithmic grids, you need an accurate solution.

- ▶ Solving large-scale problems with iterative solvers adds regularization.
- ▶ Often, scan over  $\alpha$  is done with low accuracy only.



Many other works considered very mildly ill-posed problems (e.g., denoising) only, and only considered single noise realizations.

- ▶ Unbiased risk estimators can be **problematic** for ill-posed problems.
- ▶ Asymptotic analysis suggests that **GSURE** is far off the real, reasonable risk function.
- ▶ In fact, **risk estimation is an asymptotically ill-posed problem itself**.
- ▶ Discrepancy principle was analysed in the same framework, and although often more conservative than SURE/GSURE, often more reliable.
  
- ▶ New risk estimators not based on Stein's method? Maybe not unbiased, i.e., regularized?
- ▶ LASSO: Asymptotic theory? Different GSURE risk more suitable (Bregman distances)?
- ▶ Non-Gaussian noise models?



**L, Proksch, Brune, Bissantz, Burger, Dette & Wübbeling, 2017.** *Risk Estimators for Choosing Regularization Parameters in Ill-Posed Problems - Properties and Limitations*, *submitted*, [arXiv:1701.04970](https://arxiv.org/abs/1701.04970).