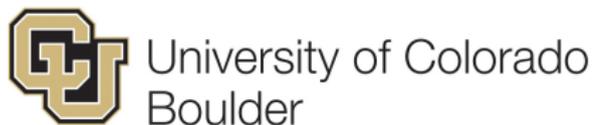


Design Matrix Uncertainty: Robust Optimization and Approximate MLE Approaches

Richard J. Clancy and Stephen Becker

University of Colorado at Boulder
Department of Applied Mathematics

July 21st
SIAM OPT 2021



Problem setup

Robust least squares

Approximate maximum likelihood estimation

Table of Contents

Problem setup

Robust least squares

Approximate maximum likelihood estimation

Problem Setup

We consider the generative model $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\eta}$

- ▶ $\mathbf{A} \in \mathbb{R}^{m \times n}$, design/data matrix
- ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$, additive noise
- ▶ $\mathbf{x} \in \mathbb{R}^n$, model parameters
- ▶ $\mathbf{y} \in \mathbb{R}^m$, measurements

Given measurements \mathbf{y} and design \mathbf{A} , estimate \mathbf{x} . Consider over-determined case where $m > n$.

Example: Estimating home prices

Generative model $\mathbf{y} = \mathbf{Ax} + \eta$

Example

- ▶ $\mathbf{A} \in \mathbb{R}^{m \times 2}$
 - Column 1: Home size (1366 sqft)
 - Column 2: Lot size (5036 sqft)
- ▶ $\mathbf{y} \in \mathbb{R}^m$: selling price for corresponding home (\$207k)
- ▶ Each row is data for a particular home
- ▶ GOAL: Estimate parameters \mathbf{x} to predict home price given info about home and lot size.

Ordinary least squares

Mismatch between modeled and true data given by $\boldsymbol{\eta} = \mathbf{A}\mathbf{x} - \mathbf{y}$.

Ordinary least squares (OLS): minimize residual sum of squares

$$\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2.$$

Closed form solution of $\mathbf{x}_{\text{OLS}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$.

Implicitly assumes Gaussian noise.

Uncertain design matrix

CONCERN: Assumes \mathbf{A} is known precisely!!!

In practice, it is uncommon to know \mathbf{A} with certainty. Causes are

- ▶ precision limits in measurement (dynamic range of sensor)
- ▶ truncation for memory savings (fixed number of sig. figs.)
- ▶ subjective features (survey responses)
- ▶ modeling error (features are approximate functions of data)

Example: Estimating home prices with uncertain data

Given the generative model $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\eta}$, where \mathbf{A} and $\boldsymbol{\eta}$ are both uncertain, estimate \mathbf{x} .

Example

- ▶ $\mathbf{A} \in \mathbb{R}^{m \times 2}$,
 - Column 1: *Rounded* home size (~~1366~~ 1400 sqft)
 - Column 2: *Rounded* lot size (~~5036~~ 5000 sqft).
- ▶ $\mathbf{y} \in \mathbb{R}^m$ selling price for corresponding home (\$207k)
- ▶ Each row is data for a particular home
- ▶ GOAL: Estimate parameters \mathbf{x} to predict home price given **uncertain** info about home and lot size.

Total least squares

Total least squares considers uncertainty in \mathbf{A} as well and solves

$$\begin{aligned} & \min_{\mathbf{U}, \boldsymbol{\eta}, \mathbf{x}} && \|\mathbf{U}, \boldsymbol{\eta}\|_F \\ & \text{subject to} && (\mathbf{A} + \mathbf{U})\mathbf{x} = \mathbf{y} + \boldsymbol{\eta}. \end{aligned}$$

- ▶ Closed form solution of $\mathbf{x}_{\text{TLS}} = (\mathbf{A}^T \mathbf{A} - s_{n+1}^2 \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$ where s_{n+1} is the smallest singular value of $[\mathbf{A}, \mathbf{y}] \in \mathbb{R}^{m \times (n+1)}$ [H. Golub and F. van Loan, 1980].
- ▶ TLS is *always* worse conditioned than OLS (which is known to be poorly conditioned)
- ▶ Often performs poorly in practice...need better method

Our contribution

To address limitations of existing methods for regression with uncertain design matrices, we present two methods:

1. using box-constrained robust optimization, and
2. an approximate MLE framework based on the saddle point approximation from complex analysis.

Both approaches allow us to move away from implicit/explicit assumptions of Gaussianity and provides data scientists with more tools for handling model uncertainty.

Table of Contents

Problem setup

Robust least squares

Approximate maximum likelihood estimation

Robust optimization

Robust optimization (RO) seeks optima over all realizations of uncertain data

$$\min_{\mathbf{x}} \max_{\mathbf{U} \in \mathcal{U}} f(\mathbf{x}, \mathbf{U})$$

where \mathcal{U} is an uncertainty set and f is a real-valued function.

- ▶ Noise comes from a bounded set
- ▶ Uncertainty has geometric rather than distributional interpretation (can be viewed as uniform noise over set)
- ▶ Considers “worst-case” data

Robust least squares

Let $f(\mathbf{x}, \mathbf{U}) = \|(\mathbf{A} + \mathbf{U})\mathbf{x} - \mathbf{y}\|^2$ and solve the robust least squares problem

$$\min_{\mathbf{x}} \left\{ \underbrace{\max_{\|\mathbf{U}\|_{\infty} \leq \delta} \|(\mathbf{A} + \mathbf{U})\mathbf{x} - \mathbf{y}\|^2}_{F(\mathbf{x})} \right\}$$

where $\|\mathbf{U}\|_{\infty}$ is the largest magnitude element of matrix \mathbf{U} .
Parameter δ can be inferred from observed data.

Other works in robust least squares:

- ▶ constrained Frobenius norm [El Ghaoui and Lebret, 1997],
- ▶ constrained 1-norm [Xu et al., 2010],
- ▶ and more [Bertsimas et al., 2011].

Reformulate objective

Must optimize

$$\min_{\mathbf{x}} \left\{ \underbrace{\max_{\|\mathbf{U}\|_{\infty} \leq \delta} \|(\mathbf{A} + \mathbf{U})\mathbf{x} - \mathbf{y}\|^2}_{F(\mathbf{x})} \right\}.$$

Theorem

The inner maximization can be written as

$$F(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + 2\delta\|\mathbf{x}\|_1\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_1 + m\delta^2\|\mathbf{x}\|_1^2.$$

with the argmax given by $\mathbf{U}_{\mathbf{x}}^ = \delta \cdot \text{sign}(\mathbf{x}(\mathbf{A}\mathbf{x} - \mathbf{y})^T)$ for a fixed \mathbf{x} .*

Rewrite robust least squares as

$$\min_{\mathbf{x}} \underbrace{\{\|\mathbf{Ax} - \mathbf{y}\|^2 + 2\delta\|\mathbf{x}\|_1 \|\mathbf{Ax} - \mathbf{y}\|_1 + m\delta^2\|\mathbf{x}\|_1^2\}}_{F(\mathbf{x})}.$$

Not differentiable, but convex since max of convex functions.

Lemma

Element of subdifferential given by

$$F'(\mathbf{x}) = 2(\mathbf{A} + \mathbf{U}_x^*)^T [(\mathbf{A} + \mathbf{U}_x^*)\mathbf{x} - \mathbf{y}] \in \partial F(\mathbf{x}).$$

Access to subdifferential \rightarrow tractable

For non-smooth problems, appropriate to use:

- ▶ Subgradient method
- ▶ Bundle method

In practice, slow to converge.

Quasi-Newton methods performed well with considerable speed-up and were used for numerical experiments.

Numerical Experiments

$$\min_{\mathbf{x}} \left\{ \|\mathbf{Ax} - \mathbf{y}\|^2 + 2\delta\|\mathbf{x}\|_1 \|\mathbf{Ax} - \mathbf{y}\|_1 + m\delta^2\|\mathbf{x}\|_1^2 \right\}$$

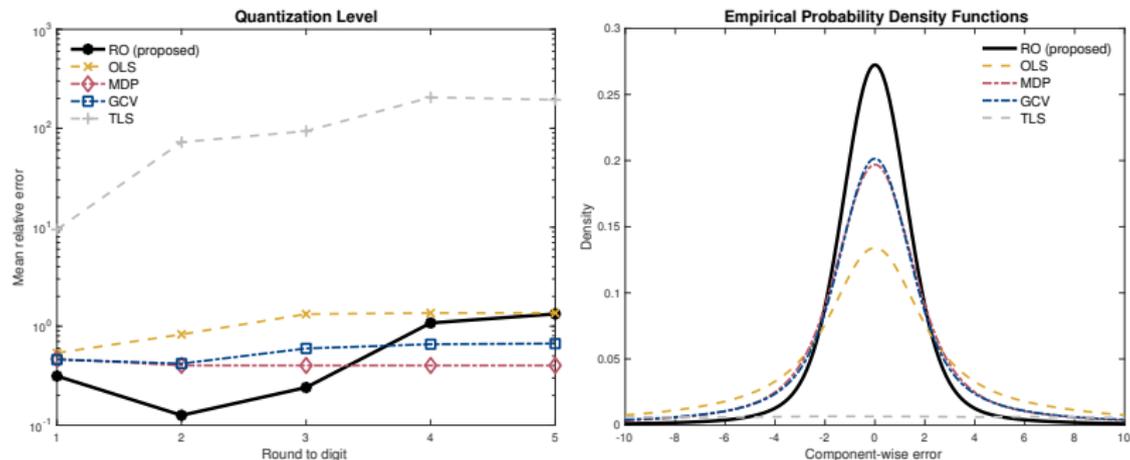


Figure: Left: mean relative error $e = \|\mathbf{x} - \hat{\mathbf{x}}\|/\|\mathbf{x}\|$ over 10k simulations for different estimators as a function of quantization level. Right: empirical PDF for errors of different methods.

Table of Contents

Problem setup

Robust least squares

Approximate maximum likelihood estimation

Maximum likelihood estimation

Once again, start with the generative model

$$\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$$

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, random/uncertain design matrix
- ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$, additive noise
- ▶ $\mathbf{x} \in \mathbb{R}^n$, model parameters
- ▶ $\mathbf{y} \in \mathbb{R}^m$, measurements

Unlike robust LS, we assume distributional knowledge of \mathbf{G} and $\boldsymbol{\eta}$, use to find PDF of \mathbf{y}

Maximum likelihood estimation

Maximum likelihood estimation (MLE) works as follows:

- ▶ Observe measurements \mathbf{y} from some distribution $\mathcal{P}_{\mathbf{y}}$ determined by distributions $\mathcal{P}_{\mathbf{G}}$ and $\mathcal{P}_{\boldsymbol{\eta}}$
- ▶ Form likelihood function from PDF of \mathbf{y} , i.e. $L(\mathbf{x}) = f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$, where $f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$ is the PDF of \mathbf{y}
- ▶ What parameters \mathbf{x} best explain observed data \mathbf{y} ?
- ▶ Find out by maximizing the likelihood function

$$\operatorname{argmax}_{\mathbf{x}} \{L(\mathbf{x}) = f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})\}$$

Maximum likelihood estimation

- ▶ We assume independence of G_{ij} 's and η_i 's throughout, implies independence of y_i 's
- ▶ When components of \mathbf{y} are independent, we can split PDF

$$f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x}) = \prod_{i=1}^m f_{Y_i}(y_i; \mathbf{x})$$

- ▶ Focus on maximizing the log-likelihood function

$$\hat{\mathbf{x}}_{\text{MLE}} = \underset{\mathbf{x}}{\operatorname{argmax}} \left\{ \sum_{i=1}^m \ln [f_{Y_i}(y_i; \mathbf{x})] \right\}$$

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{A} is known precisely

- ▶ $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$, negative log-likelihood for Gaussian noise (OLS)
- ▶ $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_1$, negative log-likelihood for Laplacian noise
- ▶ $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_\infty$, negative log-likelihood for uniform noise
- ▶ For uncertainty in operator, solving TLS problem

$$\min_{\mathbf{x}, \mathbf{U}, \boldsymbol{\eta}} \quad \|\mathbf{U}, \boldsymbol{\eta}\|_F$$

Subject to $(\mathbf{A} + \mathbf{U})\mathbf{x} = \mathbf{y} + \boldsymbol{\eta}$

yields MLE for i.i.d. Gaussian in \mathbf{U} and $\boldsymbol{\eta}$

MLE for uncertainty in design matrix

- ▶ GOAL: find PDF for $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$ to form a likelihood function
- ▶ Each component can be rewritten as a sum of random variables, i.e., $y_i = \mathbf{g}_i^T \mathbf{x} + \eta_i = \sum_{j=1}^n G_{ij}x_j + \eta_i$
- ▶ PDFs for sums of random variables are challenging to derive since they require convolutions

Example

Let $Z \sim \mathcal{N}(0, 1)$ and $U \sim \text{Uniform}(0, 1)$. The PDF for $U + Z$ is

$$f_{U+Z}(t) = \frac{1}{\sqrt{2\pi}} \int_0^1 e^{-(t-s)^2/2} ds.$$

No analytic form!

Moment generating functions

- ▶ The moment generating function (MGF) is a bilateral Laplace transform of PDF

$$M_Y(t) = \mathbb{E}(e^{tY})$$

- ▶ When MGF exists, it uniquely characterizes the distribution
- ▶ Useful properties: for independent random variables U and Z and $a \in \mathbb{R}$

$$M_{aU+Z}(t) = M_U(at)M_Z(t),$$

with MGFs, convolution \rightarrow multiplication

Example

Let $Z \sim \mathcal{N}(0, 1)$ and $U \sim \text{Uniform}(0, 1)$. The MGF for $U + Z$ is

$$M_{U+Z}(t) = \frac{(e^t - 1)e^{-t^2/2}}{t}.$$

No need for quadrature

Moment generating functions

- ▶ Using properties of MGFs,

$$M_{Y_i}(t) = M_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = M_{\eta_i}(t) \prod_{j=1}^n M_{G_{ij}}(tx_j)$$

Express complicated MGF for Y_i as the product of univariate MGFs for G_{ij} and η_i .

- ▶ By inverting transform for M_{Y_i} , we can recover density, but difficult in practice
- ▶ Approximate PDF instead

Density approximation

Options to approximate PDFs:

- ▶ Kernel density estimation : data intensive
- ▶ Edgeworth series: poor tail behavior (poly. series)
- ▶ Saddle point approximation [Daniels, 1954, Strawderman et al., 1996]: uses complex analysis, works well in practice

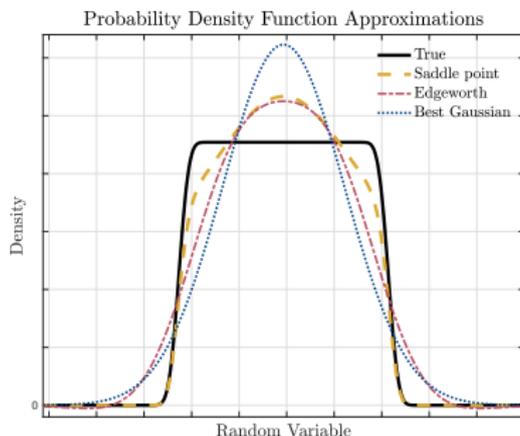


Figure: True density and several approximations for $y = \mathbf{g}^T \mathbf{x} + \eta$

Saddle point approximation

Saddle point approximation for PDF of Y is

$$f_Y(y) \approx \sqrt{\frac{1}{2\pi K_Y''(t_0)}} e^{K_Y(t_0) - yt_0}$$

- ▶ $K_Y(t) = \ln M_Y(t)$ is *Cumulant Generating Function* (CGF)
- ▶ t_0 is the solution to $K_Y'(t) - y = 0$ (use Newton's method)

Approximate MLE

Using the saddle point approximation and eliminating constants, write the approximate log-likelihood function, $\ell(\mathbf{x}) \approx \ln L(\mathbf{x})$, as

$$\begin{aligned}\ell(\mathbf{x}) &= \sum_{i=1}^m \ln \underbrace{\left\{ \sqrt{\frac{1}{K''_{Y_i}(t_i)}} e^{K_{Y_i}(t_i) - y_i t_i} \right\}}_{\approx \text{constant} \cdot f_{Y_i}(y_i)} \\ &= \sum_{i=1}^m \left\{ K_{Y_i}(t_i) - t_i y_i - \frac{1}{2} \ln (K''_{Y_i}(t_i)) \right\} \\ &= \sum_{i=1}^m \left[K_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t_i(\mathbf{x})) - \frac{1}{2} \ln (K''_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t_i(\mathbf{x}))) - t_i(\mathbf{x}) y_i \right].\end{aligned}$$

where \mathbf{g}_i^T is i^{th} row of \mathbf{G} and $t_i(\mathbf{x})$ is solution to $K'_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = y_i$.

Optimization problem

The approximate MLE can be cast generically in vector form

$$\begin{aligned} \underset{\mathbf{x}, \mathbf{t}}{\operatorname{argmax}} \quad & \mathbf{1}^T \left(K_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t}) - \frac{1}{2} \ln (K''_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t})) \right) - \mathbf{t}^T \mathbf{y} \\ \text{subject to} \quad & K'_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t}) = \mathbf{y} \end{aligned}$$

Example

When $\mathbf{G} \sim \text{Uniform}(\mathbf{H} - \delta \mathbf{1}\mathbf{1}^T, \mathbf{H} + \delta \mathbf{1}\mathbf{1}^T)$ and $\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$,

$$\begin{aligned} \underset{\mathbf{x}, \mathbf{t}}{\operatorname{argmax}} \quad & \mathbf{t}^T \left(\frac{\sigma^2}{2} \mathbf{t} + \mathbf{H} \mathbf{x} - \mathbf{y} \right) + \mathbf{1}^T \ln [\sinh (\delta \mathbf{t} \mathbf{x}^T) \odot (\delta \mathbf{t} \mathbf{x}^T)] \mathbf{1} \\ & - \frac{1}{2} \mathbf{1}^T \ln [\sigma^2 \mathbf{1} - \delta^2 \operatorname{csch}^2 (\delta \mathbf{t} \mathbf{x}^T) \mathbf{x}^2] \\ \text{Subject to} \quad & \sigma^2 \mathbf{t} + \mathbf{H} \mathbf{x} + \delta \operatorname{coth} (\delta \mathbf{t} \mathbf{x}^T) \mathbf{x} - n(\mathbf{1} \odot \mathbf{t}) = \mathbf{y}. \end{aligned}$$

Gradients

Despite constraint requiring a numeric solve, have gradients.

Letting $\mathbf{q}(\mathbf{x}, \mathbf{t}) = K'_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) - \mathbf{y}$ be our constraint, using the chain rule and implicit differentiation we have

$$\nabla_{\mathbf{x}} \ell = \frac{\partial \ell}{\partial \mathbf{x}} - \left(\frac{\partial \ell}{\partial \mathbf{t}} \right) \left(\frac{\partial \mathbf{q}}{\partial \mathbf{t}} \right)^{-1} \left(\frac{\partial \mathbf{q}}{\partial \mathbf{x}} \right).$$

Unimportant, but for completeness, each factor is given by:

$$\frac{\partial \ell}{\partial \mathbf{x}} = \mathbf{1}^T \left(\frac{\partial}{\partial \mathbf{x}} K_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) - \frac{1}{2} \{ \text{diag} (K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t})) \}^{-1} \frac{\partial}{\partial \mathbf{x}} K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) \right)$$

$$\frac{\partial \ell}{\partial \mathbf{t}} = \left(K'_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) - \frac{1}{2} (K'''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) \oslash K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t})) - \mathbf{y} \right)^T$$

$$\frac{\partial \mathbf{q}}{\partial \mathbf{t}} = \text{diag} (K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t})) ,$$

$$\frac{\partial \mathbf{q}}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{x}} K'_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}).$$

Numerical experiments

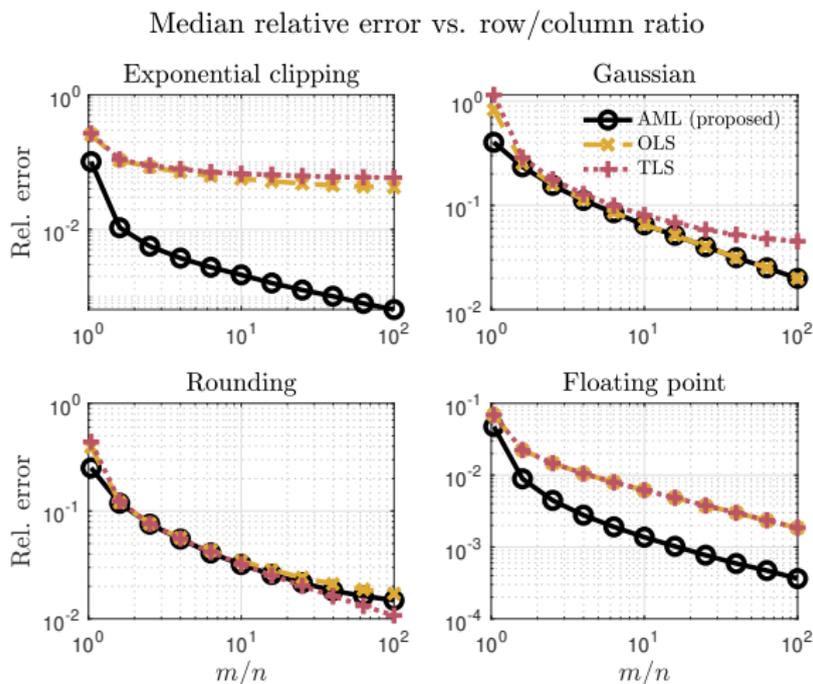


Figure: Median relative error over 1,000 simulations as number of rows increase for fixed number of columns in design matrix.

Numerical Experiments

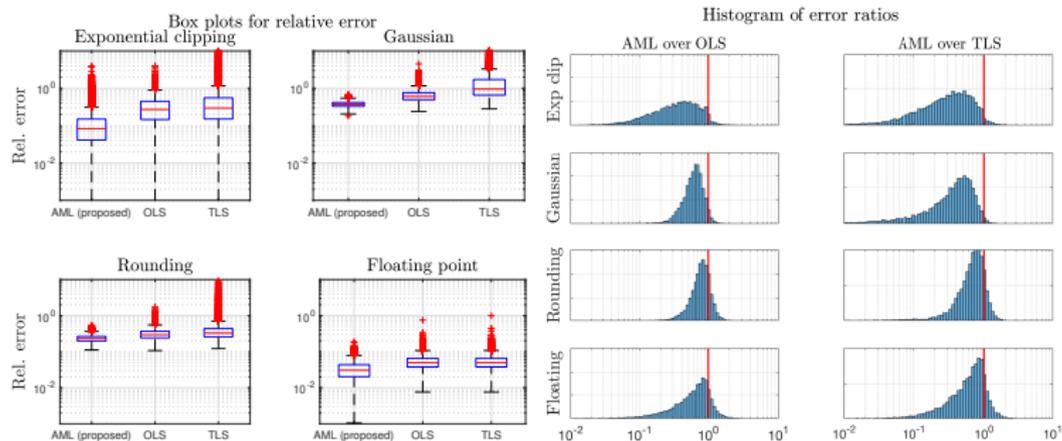


Figure: Error metrics for simulations $\mathbf{G} \in \mathbb{R}^{55 \times 50}$. Left: box-plot of relative error for different methods. Right: histogram of error ratio $\|\mathbf{x}_{\text{AML}} - \mathbf{x}_{\text{TRU}}\| / \|\mathbf{x}_{\text{OLS}} - \mathbf{x}_{\text{TRU}}\|$. Values less than 1 indicate AML outperformed competing method for identical data.

Summary

- ▶ Motivated the need for solving regression problems with uncertain design matrices.
- ▶ Formulated and solved a robust least squares problem that is well suited for regression problems with data subject to quantization error.
- ▶ Derived an approximate MLE function and provided its gradient. Ideal for regression problems with different noise models across features.

References

- D. Bertsimas, D. B. Brown, and C. Caramanis. Theory and applications of robust optimization. SIAM Review, 53 (3):464–501, 2011. doi: 10.1137/080734510. URL <https://doi.org/10.1137/080734510>.
- Henry E Daniels. Saddlepoint approximations in statistics. The annals of mathematical statistics, pages 631–650, 1954.
- Laurent El Ghaoui and Herve Le Bret. Robust solutions to least-squares problems with uncertain data. SIAM J. Matrix Anal. Appl., 18(4):1035 – 1064, 10 1997.
- Gene H. Golub and Charles F. van Loan. An analysis of the total least squares problem. SIAM J. Numer. Anal., 17 (6):883 – 893, 12 1980.
- Robert L Strawderman, George Casella, and Martin T Wells. Practical small-sample asymptotics for regression problems. Journal of the American Statistical Association, 91(434):643–654, 1996.
- H. Xu, C. Caramanis, and S. Mannor. Robust regression and lasso. IEEE Tran. Info. Theory, 56(7):3561–3574, July 2010. ISSN 0018-9448. doi: 10.1109/TIT.2010.2048503.

Thank you for your time!

- ▶ Stephen Becker, and Richard J. Clancy. “Robust least squares for quantized data matrices.” *Signal Processing* 176 (2020)
- ▶ Richard J. Clancy and Stephen Becker. “Approximate maximum likelihood estimators for linear regression with design matrix uncertainty.” *arXiv preprint arXiv:2104.03307* (2021).

Questions?