

Stein unbiased risk estimators
for tuning hyperparameters
of distributed regression algorithms

Damiano Varagnolo, ITK



Gianluigi Pillonetto



Luca Schenato

Purposes of this seminar

- 1 discuss about a useful tool
- 2 connect with you

Roadmap

- Stein's lemma
- Stein's unbiased risk estimator (SURE)
- Model selection through SURE
(will follow "Stein's Unbiased Risk Estimate, Statistical Machine Learning 2015, Tibshirani & Wasserman")
- RKHS-based regression
- average-consensus algorithms
- SURE in our distributed regression context

Stein's lemma

Stein's univariate lemma

If:

$$X \sim \mathcal{N}(\mu, \sigma^2) \quad (1)$$

$$f : \mathbb{R} \mapsto \mathbb{R} \text{ absolutely continuous} \quad (2)$$

$$f' \text{ exists and is s.t. } \mathbb{E}[|f'(X)|] < +\infty \quad (3)$$

Stein's univariate lemma

If:

$$X \sim \mathcal{N}(\mu, \sigma^2) \quad (1)$$

$$f : \mathbb{R} \mapsto \mathbb{R} \text{ absolutely continuous} \quad (2)$$

$$f' \text{ exists and is s.t. } \mathbb{E}[|f'(X)|] < +\infty \quad (3)$$

then:

$$\mathbb{E}[(X - \mu) f(X)] = \sigma^2 \mathbb{E}[f'(X)] \quad (4)$$

Implications

(with $\sigma^2 = 1$ for notational simplicity)

$$\mathbb{E}[(X - \mu) f(X)] = \text{cov}(X, f(X)) = \mathbb{E}[f'(X)] \quad (5)$$

Implications

(with $\sigma^2 = 1$ for notational simplicity)

$$\mathbb{E}[(X - \mu) f(X)] = \text{cov}(X, f(X)) = \mathbb{E}[f'(X)] \quad (5)$$

- if f is an estimator of μ , then estimating $\text{cov}(X, f(X))$ through $\mathbb{E}[(X - \mu) f(X)]$ requires knowing $\mu \implies$ not a feasible path!

Implications

(with $\sigma^2 = 1$ for notational simplicity)

$$\mathbb{E}[(X - \mu) f(X)] = \text{cov}(X, f(X)) = \mathbb{E}[f'(X)] \quad (5)$$

- if f is an estimator of μ , then estimating $\text{cov}(X, f(X))$ through $\mathbb{E}[(X - \mu) f(X)]$ requires knowing $\mu \implies$ not a feasible path!
- alternative strategy: estimate $\text{cov}(X, f(X))$ through computing $\widehat{\mathbb{E}}[f'(X)]$

Stein's multivariate lemma

If:

$$X \sim \mathcal{N}(\mu, \sigma^2 I) \quad (6)$$

$f : \mathbb{R}^n \mapsto \mathbb{R}$ *almost differentiable*, i.e.:

$f(\cdot, x_{-i}) : \mathbb{R} \mapsto \mathbb{R}$ absolutely continuous for a.e. $x_i \in \mathbb{R}^{n-1}$ and $i = 1, \dots, n$ (7)

$$\frac{\partial f}{\partial x_i} \text{ exists and is s.t. } \mathbb{E} \left[\left| \frac{\partial f}{\partial x_i}(X) \right| \right] < +\infty \text{ for each } i = 1, \dots, n \quad (8)$$

Stein's multivariate lemma

If:

$$X \sim \mathcal{N}(\mu, \sigma^2 I) \quad (6)$$

$f : \mathbb{R}^n \mapsto \mathbb{R}$ *almost differentiable*, i.e.:

$f(\cdot, x_{-i}) : \mathbb{R} \mapsto \mathbb{R}$ absolutely continuous for a.e. $x_i \in \mathbb{R}^{n-1}$ and $i = 1, \dots, n$ (7)

$$\frac{\partial f}{\partial x_i} \text{ exists and is s.t. } \mathbb{E} \left[\left| \frac{\partial f}{\partial x_i}(X) \right| \right] < +\infty \text{ for each } i = 1, \dots, n \quad (8)$$

then:

$$\mathbb{E}[(X - \mu) f(X)] = \sigma^2 \mathbb{E}[\nabla f(X)] \quad (9)$$

Stein's multivariate lemma (2)

$$X \sim \mathcal{N}(\mu, \sigma^2 I) \quad f: \mathbb{R}^n \mapsto \mathbb{R}^n \quad f = [f_1, \dots, f_n] \quad (10)$$

$$\implies \mathbb{E}[(X - \mu) f_i(X)] = \sigma^2 \mathbb{E}[\nabla f_i(X)] \quad (11)$$

Stein's multivariate lemma (2)

$$X \sim \mathcal{N}(\mu, \sigma^2 I) \quad f : \mathbb{R}^n \mapsto \mathbb{R}^n \quad f = [f_1, \dots, f_n] \quad (10)$$

$$\implies \mathbb{E}[(X - \mu) f_i(X)] = \sigma^2 \mathbb{E}[\nabla f_i(X)] \quad (11)$$

$$\implies \sum_{i=1}^n \text{cov}(X_i, f_i(X)) = \sigma^2 \mathbb{E} \left[\sum_{i=1}^n \frac{\partial f_i}{\partial x_i}(X) \right] \quad (12)$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \hat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \hat{\mu}(y) = \text{estimate of } \mu \text{ at } y \quad (13)$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \hat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \hat{\mu}(y) = \text{estimate of } \mu \text{ at } y \quad (13)$$

$$R = \mathbb{E} [\|\mu - \hat{\mu}\|_2^2]$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \hat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \hat{\mu}(y) = \text{estimate of } \mu \text{ at } y \quad (13)$$

$$\begin{aligned} R &= \mathbb{E} [\|\mu - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y + y - \hat{\mu}\|_2^2] \end{aligned}$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \hat{\mu}: \mathbb{R}^n \mapsto \mathbb{R}^n \quad \hat{\mu}(y) = \text{estimate of } \mu \text{ at } y \quad (13)$$

$$\begin{aligned} R &= \mathbb{E} [\|\mu - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y + y - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y\|_2^2] + \mathbb{E} [\|y - \hat{\mu}\|_2^2] + 2\mathbb{E} [(\mu - y)^T (y - \hat{\mu})] \end{aligned}$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \hat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \hat{\mu}(y) = \text{estimate of } \mu \text{ at } y \quad (13)$$

$$\begin{aligned} R &= \mathbb{E} [\|\mu - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y + y - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y\|_2^2] + \mathbb{E} [\|y - \hat{\mu}\|_2^2] + 2\mathbb{E} [(\mu - y)^T (y - \hat{\mu})] \\ &= n\sigma^2 + \mathbb{E} [\|y - \hat{\mu}\|_2^2] + 2\mathbb{E} [(\mu - y)^T (y - \hat{\mu})] \end{aligned}$$

From Stein's multivariate lemma to Stein's unbiased risk estimate (SURE)

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \hat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \hat{\mu}(y) = \text{estimate of } \mu \text{ at } y \quad (13)$$

$$\begin{aligned} R &= \mathbb{E} [\|\mu - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y + y - \hat{\mu}\|_2^2] \\ &= \mathbb{E} [\|\mu - y\|_2^2] + \mathbb{E} [\|y - \hat{\mu}\|_2^2] + 2\mathbb{E} [(\mu - y)^T (y - \hat{\mu})] \\ &= n\sigma^2 + \mathbb{E} [\|y - \hat{\mu}\|_2^2] + 2\mathbb{E} [(\mu - y)^T (y - \hat{\mu})] \\ &= -n\sigma^2 + \mathbb{E} [\|y - \hat{\mu}\|_2^2] + 2 \sum_{i=1}^n \text{cov}(y_i, \hat{\mu}_i) \end{aligned} \quad (14)$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$R = \mathbb{E} [\|\mu - \widehat{\mu}\|_2^2] = -n\sigma^2 + \mathbb{E} [\|y - \widehat{\mu}\|_2^2] + 2 \sum_{i=1}^n \text{cov} (y_i, \widehat{\mu}_i)$$

From Stein's multivariate lemma
to Stein's unbiased risk estimate (SURE)

$$R = \mathbb{E} [\|\mu - \widehat{\mu}\|_2^2] = -n\sigma^2 + \mathbb{E} [\|y - \widehat{\mu}\|_2^2] + 2 \sum_{i=1}^n \text{cov} (y_i, \widehat{\mu}_i) \quad (15)$$
$$\implies \widehat{R} = -n\sigma^2 + \|y - \widehat{\mu}\|_2^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{\mu}_i}{\partial y_y} (y)$$

$$\text{with } \mathbb{E} [\widehat{R}] = R \quad (16)$$

Model selection through SURE in general

Model selection through SURE in general

$$\widehat{\mu} \mapsto \widehat{\mu}_\lambda \quad (17)$$

Model selection through SURE in general

$$\widehat{\mu} \mapsto \widehat{\mu}_\lambda \quad (17)$$

$$\widehat{R}_\lambda = -n\sigma^2 + \|y - \widehat{\mu}_\lambda\|_2^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{\mu}_{\lambda,i}}{\partial y_i}(y) \quad (18)$$

Model selection through SURE in general

$$\widehat{\mu} \mapsto \widehat{\mu}_\lambda \quad (17)$$

$$\widehat{R}_\lambda = -n\sigma^2 + \|y - \widehat{\mu}_\lambda\|_2^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{\mu}_{\lambda,i}}{\partial y_i}(y) \quad (18)$$

$$\lambda^* = \arg \min_{\lambda \in \Lambda} \|y - \widehat{\mu}_\lambda\|_2^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{\mu}_{\lambda,i}}{\partial y_i}(y) \quad (19)$$

requires:

- to verify that $\widehat{\mu}_\lambda$ is almost differentiable

Model selection through SURE in general

$$\widehat{\mu} \mapsto \widehat{\mu}_\lambda \quad (17)$$

$$\widehat{R}_\lambda = -n\sigma^2 + \|y - \widehat{\mu}_\lambda\|_2^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{\mu}_{\lambda,i}}{\partial y_i}(y) \quad (18)$$

$$\lambda^* = \arg \min_{\lambda \in \Lambda} \|y - \widehat{\mu}_\lambda\|_2^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{\mu}_{\lambda,i}}{\partial y_i}(y) \quad (19)$$

requires:

- to verify that $\widehat{\mu}_\lambda$ is almost differentiable
- to compute the *divergence* of $\widehat{\mu}_\lambda$, i.e., $\frac{\partial \widehat{\mu}_{\lambda,i}}{\partial y_i}(y)$

Model selection through SURE: the linear case

If:

$$y = \mu + e \quad \mu \text{ deterministic} \quad (20)$$

$$y^* = \mu + e^* \text{ future measurements on the same input locations} \quad (21)$$

$$e, e^* \text{ uncorrelated, zero mean, with covariance } \Sigma \quad (22)$$

$$\hat{y} = Sy \text{ linear estimator of } y^* \quad (23)$$

Model selection through SURE: the linear case

If:

$$y = \mu + e \quad \mu \text{ deterministic} \quad (20)$$

$$y^* = \mu + e^* \text{ future measurements on the same input locations} \quad (21)$$

$$e, e^* \text{ uncorrelated, zero mean, with covariance } \Sigma \quad (22)$$

$$\widehat{y} = Sy \text{ linear estimator of } y^* \quad (23)$$

Then:

$$\|y - \widehat{y}\|^2 + 2\text{tr}(S\Sigma) \text{ is an unbiased estimator of the risk } \mathbb{E} \left[\|y^* - \widehat{y}\|_2^2 \right] \quad (24)$$

Model selection through SURE: examples of literature

literature on other cases:

- Li (1985), *From Stein's unbiased risk estimates to the method of generalized cross-validation*, Annals of Statistics
- Li (1986), *Asymptotic optimality of c_l and generalized cross-validation in ridge regression with application to spline smoothing*, Annals of Statistics
- Johnstone (1986), *On inadmissibility of some unbiased estimates of loss*, technical report
- Kneip (1994), *Ordered linear smoothers*, Annals of Statistics
- Donoho & Johnstone (1995), *Adapting to unknown smoothness via wavelet shrinkage*, Journal of the American Statistical Association
- Efron et al. (2004), *Least angle regression*, Annals of Statistics
- Zou et al. (2007), *On the degrees of freedom of the lasso*, Annals of Statistics
- Tibshirani & Taylor (2011), *The solution path of the generalized lasso*, Annals of Statistics
- Tibshirani & Taylor (2012), *Degrees of freedom in lasso problems*, Annals of Statistics

Model selection through SURE for a specific case

The original practical problem: function estimation

$$f : \mathcal{X} \rightarrow \mathbb{R} \quad (25)$$

$$y_m = f(x_m) + \nu_m \quad m = 1, \dots, M \quad (26)$$

$$x_m \sim \mu(\mathcal{X}) \text{ i.i.d.} \quad \nu_m \sim \mathcal{N}(0, \sigma_\nu^2) \quad m = 1, \dots, M \quad (27)$$

$$\{x_m\}_{m=1}^M \quad \{\nu_m\}_{m=1}^M \quad \text{mutually independent} \quad (28)$$

The original practical problem: function estimation

$$f : \mathcal{X} \rightarrow \mathbb{R} \quad (25)$$

$$y_m = f(x_m) + \nu_m \quad m = 1, \dots, M \quad (26)$$

$$x_m \sim \mu(\mathcal{X}) \text{ i.i.d.} \quad \nu_m \sim \mathcal{N}(0, \sigma_\nu^2) \quad m = 1, \dots, M \quad (27)$$

$$\{x_m\}_{m=1}^M \quad \{\nu_m\}_{m=1}^M \quad \text{mutually independent} \quad (28)$$

Problem: estimate f starting from $\{x_m, y_m\}$

Nonparametric approach *(cast the problem as a Gaussian regression)*

$$f \sim \mathcal{N}(0, K) \quad K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad \text{so that} \quad \mathbb{E}[f(x)f(x')] = K(x, x') \quad (29)$$

Nonparametric approach *(cast the problem as a Gaussian regression)*

$$f \sim \mathcal{N}(0, K) \quad K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad \text{so that} \quad \mathbb{E}[f(x)f(x')] = K(x, x') \quad (29)$$

Examples:

- Brownian motion: $K(x, x') = \min(x, x')$ $\mathcal{X} = [0, 1]$
- Radial basis: $K(x, x') = \exp(-\|x - x'\|^2)$ $\mathcal{X} \subseteq \mathbb{R}^m$

Maximum a posteriori estimator

$$\hat{f}_{\text{MAP}}(x) = \mathbb{E} \left[f(x) \mid \{x_m, y_m\} \right] \quad (\text{also MV})$$

Maximum a posteriori estimator

$$\begin{aligned}\widehat{f}_{\text{MAP}}(x) &= \mathbb{E} \left[f(x) \mid \{x_m, y_m\} \right] && \text{(also MV)} \\ &= \sum_{m=1}^M K(x, x_m) c_m && \text{(a.k.a. regularization network)}\end{aligned}\tag{30}$$

Maximum a posteriori estimator

$$\begin{aligned}\widehat{f}_{\text{MAP}}(x) &= \mathbb{E} \left[f(x) \mid \{x_m, y_m\} \right] && \text{(also MV)} \\ &= \sum_{m=1}^M K(x, x_m) c_m && \text{(a.k.a. regularization network)}\end{aligned}\tag{30}$$

$$\begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = H_{\text{MAP}} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}\tag{31}$$

$$H_{\text{MAP}} := \left(\begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix} + \sigma_\nu^2 I \right)^{-1}\tag{32}$$

Gaussian regression – practical issues associated to the MAP

$$\widehat{f}_{\text{MAP}}(x) = [K(x, x_1) \ \dots \ K(x, x_M)] H_{\text{MAP}} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} \quad (33)$$

$$H_{\text{MAP}} := \left(\begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \dots & K(x_M, x_M) \end{bmatrix} + \sigma_\nu^2 I \right)^{-1} \quad (34)$$

computational cost $O(M^3)$

How may we tackle the $O(M^3)$ computational cost issue?

$$H_{\text{MAP}} := \left(\begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix} + \sigma_\nu^2 I \right)^{-1} \quad (35)$$

Typical approaches: low-rank / sparsification approximations



Smola & Schölkopf (2000)

Sparse greedy matrix approximations for machine learning



Quiñero-Candela & Rasmussen (2005)

A unifying view of sparse approximate Gaussian process regression



Bach & Jordan (2005)

Predictive low-rank decompositions for kernel methods



Snelson & Ghahramani (2006)

Sparse Gaussian processes using pseudo inputs



Culis et al. (2006)

Learning low-rank kernel matrices



Zhang & Kwok (2010)

Clustered Nyström method for large scale manifold learning and dimension reduction



Ambikasaran et al. (2016)

Fast direct methods for Gaussian processes

Our approach: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x)$$

Our approach: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=: \text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=: \text{remainder}} \quad (36)$$

Our approach: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=: \text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=: \text{remainder}} \quad (36)$$

$$\lambda_e \phi_e(x) = \int_{\mathcal{X}} K(x, x') \phi_e(x') d\mu(x') \quad \lambda_1 \geq \lambda_2 \dots > 0 \quad (37)$$

$$K(x, x') = \sum_{e=1}^{+\infty} \lambda_e \phi_e(x) \phi_e(x') \quad \int_{\mathcal{X}} \phi_i(x) \phi_j(x) d\mu(x) = \delta_{ij} \quad (38)$$



Our approach: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=:\text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=:\text{remainder}} \quad (36)$$

$$\lambda_e \phi_e(x) = \int_{\mathcal{X}} K(x, x') \phi_e(x') d\mu(x') \quad \lambda_1 \geq \lambda_2 \dots > 0 \quad (37)$$

$$K(x, x') = \sum_{e=1}^{+\infty} \lambda_e \phi_e(x) \phi_e(x') \quad \int_{\mathcal{X}} \phi_i(x) \phi_j(x) d\mu(x) = \delta_{ij} \quad (38)$$

$$a_e \sim \mathcal{N}(0, \lambda_e), \quad e = 1, \dots, E \quad b_e \sim \mathcal{N}(0, \lambda_{E+e}), \quad e = 1, 2, \dots \quad (39)$$



Zhu et al. (1998)

Gaussian regression and optimal finite dimensional linear models

Our approach: Karhunen-Loève expansions

$$f(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=: \text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=: \text{remainder}} \quad \lambda_e \phi_e(x) = \int_{\mathcal{X}} K(x, x') \phi_e(x') d\mu(x')$$

(40)

\implies first E ϕ_e 's = best a-priori E -dimensional approximation in a MSE sense

Our approach: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (41)$$

Our approach: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (41)$$

$$G := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix} \quad Z := \begin{bmatrix} \phi_{E+1}(x_1) & \phi_{E+2}(x_1) & \dots \\ \vdots & & \vdots \\ \phi_{E+1}(x_M) & \phi_{E+2}(x_M) & \dots \end{bmatrix} \quad (42)$$

Our approach: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (41)$$

$$G := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix} \quad Z := \begin{bmatrix} \phi_{E+1}(x_1) & \phi_{E+2}(x_1) & \dots \\ \vdots & & \vdots \\ \phi_{E+1}(x_M) & \phi_{E+2}(x_M) & \dots \end{bmatrix} \quad (42)$$

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (43)$$

Our approach: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (41)$$

$$G := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix} \quad Z := \begin{bmatrix} \phi_{E+1}(x_1) & \phi_{E+2}(x_1) & \dots \\ \vdots & & \vdots \\ \phi_{E+1}(x_M) & \phi_{E+2}(x_M) & \dots \end{bmatrix} \quad (42)$$

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (43)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \dots \quad \phi_E(x)] \widehat{\mathbf{a}} \quad \widehat{\mathbf{a}} = H\mathbf{y} \quad H := \left(\frac{G^T G}{M} + \frac{\sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (44)$$

Summary

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (45)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \cdots \quad \phi_E(x)] H \mathbf{y} \quad H := \left(\frac{G^T G}{M} + \frac{\sigma_{\nu}^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (46)$$

computational cost: $O(E^3)$

Summary

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (45)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \cdots \quad \phi_E(x)] H \mathbf{y} \quad H := \left(\frac{G^T G}{M} + \frac{\sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (46)$$

computational cost: $O(E^3)$

interesting for us because $\frac{G^T G}{M} = \frac{1}{M} \sum_{m=1}^M G_m^T G_m \quad \frac{G^T \mathbf{y}}{M} = \frac{1}{M} \sum_{m=1}^M G_m^T y_m \quad (47)$

Average consensus

(i.e., how to compute an average in a distributed fashion?)

synchronous communications

synchronous consensus: $\mathbf{x}(k+1) = P\mathbf{x}(k)$ (with P doubly stochastic) (Markov chains ('60s), Seneta 2006, ...)

synchronous communications

synchronous consensus: $\mathbf{x}(k+1) = P\mathbf{x}(k)$ (with P doubly stochastic) (Markov chains ('60s), Seneta 2006, ...)

asynchronous communications with perfect channel feedback

ratio consensus (Bénézit et al. 2010)

synchronous communications

synchronous consensus: $\mathbf{x}(k+1) = P\mathbf{x}(k)$ (with P doubly stochastic) (Markov chains ('60s), Seneta 2006, ...)

asynchronous communications with perfect channel feedback

ratio consensus (Bénézit et al. 2010)

asynchronous communications without perfect channel feedback

robust ratio consensus (Dominguez-Garcia et al. 2011)

Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

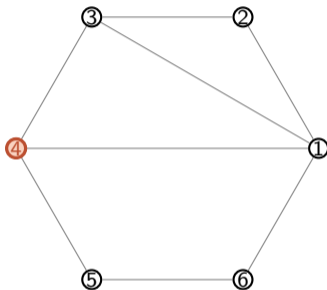
$$\left\{ \begin{array}{l} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{array} \right.$$

Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

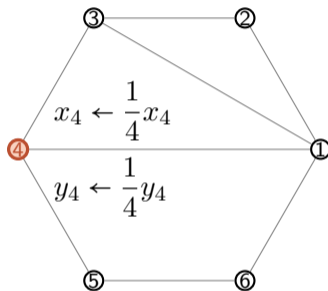


Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

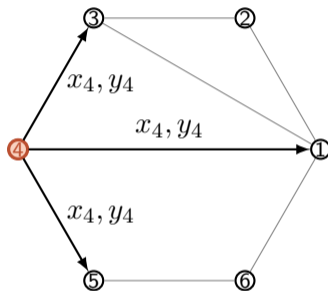


Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

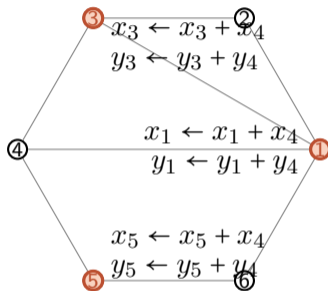


Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



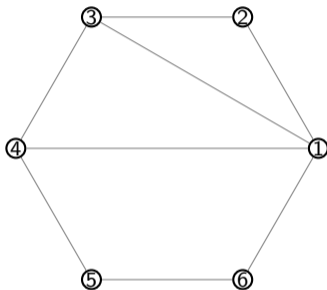
Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{cases} x_i(k) \rightarrow \beta_i(k) \sum_j x_j(0) \\ y_i(k) \rightarrow \beta_i(k) \sum_j y_j(0) \end{cases}$$

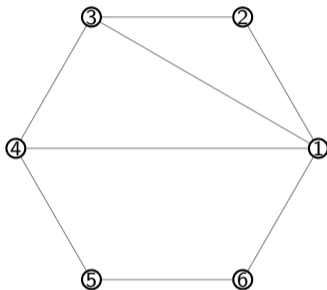


Ratio consensus

asynchronous communications with perfect channel feedback (Bénézit et al. 2010)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



$$\begin{cases} x_i(k) \rightarrow \beta_i(k) \sum_j x_j(0) \\ y_i(k) \rightarrow \beta_i(k) \sum_j y_j(0) \end{cases} \implies z_i(k) := \frac{x_i(k)}{y_i(k)} \rightarrow \frac{\sum_i x_i(0)}{\sum_i y_i(0)} = \frac{1}{N} \sum_i \theta_i$$

Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

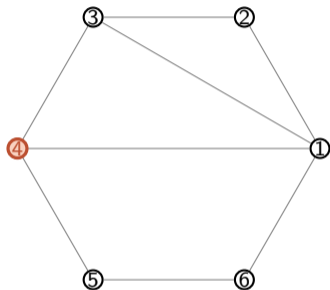
$$\left\{ \begin{array}{l} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{array} \right.$$

Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

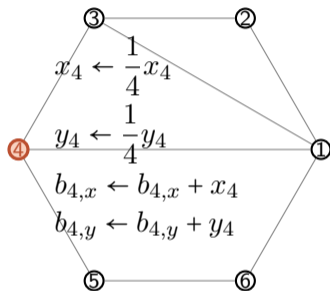
$$P(k) = \begin{bmatrix} 1 & 0 & 0 & \mathbf{0} & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$
$$P(k) = \begin{bmatrix} 1 & 0 & 0 & \mathbf{0} & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



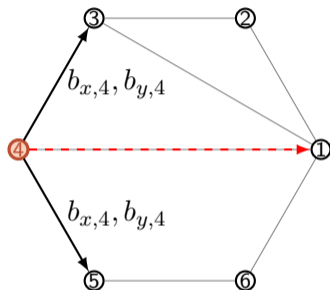
- $b_{i,x}$: total cumulative mass of x_i
- $\beta_{i,x}^{(j)}$: j 's local estimate of $b_{i,x}$

Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & \mathbf{0} & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



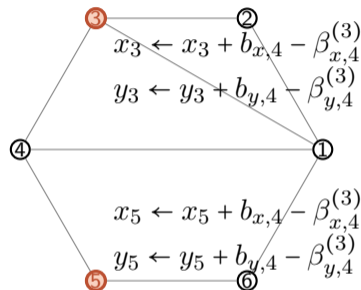
- $b_{i,x}$: total cumulative mass of x_i
- $\beta_{i,x}^{(j)}$: j 's local estimate of $b_{i,x}$

Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & \mathbf{0} & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



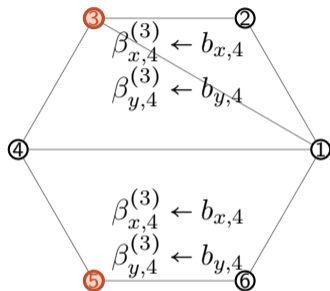
- $b_{i,x}$: total cumulative mass of x_i
- $\beta_{i,x}^{(j)}$: j 's local estimate of $b_{i,x}$

Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & \mathbf{0} & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



- $b_{i,x}$: total cumulative mass of x_i
- $\beta_{i,x}^{(j)}$: j 's local estimate of $b_{i,x}$

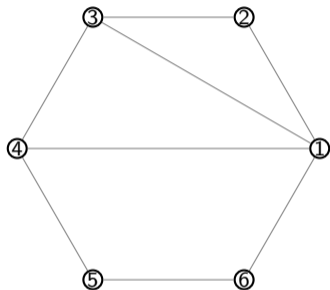
Robust ratio consensus

asynchronous communication **without** perfect channel feedback (Dominguez-Garcia et al. 2011)

$$\begin{cases} \mathbf{x}(k+1) = P(k)\mathbf{x}(k) \\ x_i(0) = \theta_i \\ \mathbf{y}(k+1) = P(k)\mathbf{y}(k) \\ y_i(0) = 1 \end{cases}$$

$$P(k) = \begin{bmatrix} 1 & 0 & 0 & \mathbf{0} & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

- $b_{i,x}$: total cumulative mass of x_i
- $\beta_{i,x}^{(j)}$: j 's local estimate of $b_{i,x}$



$$z_i(k) = \frac{x_i(k)}{y_i(k)} \rightarrow \frac{1}{N} \sum_j \theta_j$$

SURE in our distributed regression settings

Recap: KL + average consensus

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (48)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \dots \quad \phi_E(x)] H \mathbf{y} \quad H := \left(\frac{G^T G}{M} + \lambda \frac{\sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (49)$$

$$\frac{G^T G}{M} = \frac{1}{M} \sum_{m=1}^M G_m^T G_m \quad \frac{G^T \mathbf{y}}{M} = \frac{1}{M} \sum_{m=1}^M G_m^T y_m \quad (50)$$

Recap: KL + average consensus

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (48)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \dots \quad \phi_E(x)] H \mathbf{y} \quad H := \left(\frac{G^T G}{M} + \lambda \frac{\sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (49)$$

$$\frac{G^T G}{M} = \frac{1}{M} \sum_{m=1}^M G_m^T G_m \quad \frac{G^T \mathbf{y}}{M} = \frac{1}{M} \sum_{m=1}^M G_m^T y_m \quad (50)$$

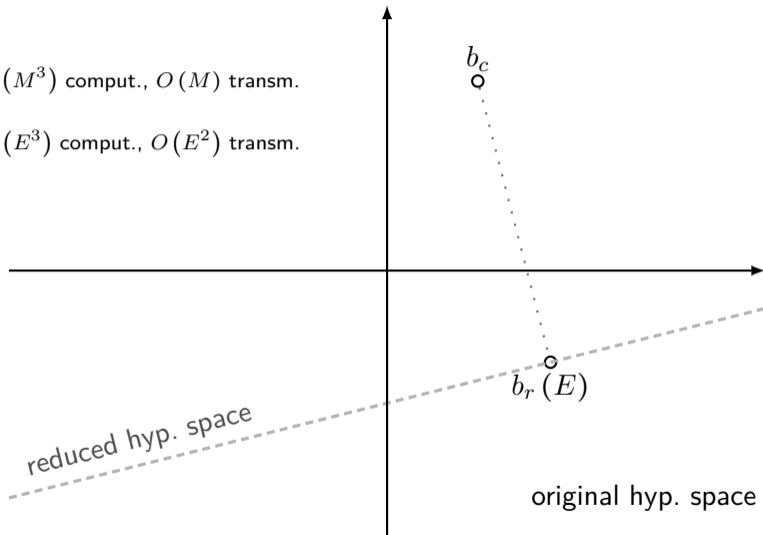
Questions:

- how shall we tune E ?
- how shall we tune λ ?

Tuning of E (not in this presentation)

b_c : $O(M^3)$ comput., $O(M)$ transm.

b_r : $O(E^3)$ comput., $O(E^2)$ transm.



Tuning of λ : a SURE-based approach

Recall:

$$\|\mathbf{y} - \widehat{\mathbf{y}}\|^2 + 2\text{tr}(S\Sigma) \text{ is an unbiased estimator of the risk } \mathbb{E}\left[\|\mathbf{y}^* - \widehat{\mathbf{y}}\|_2^2\right] \text{ with } \widehat{\mathbf{y}} = S\mathbf{y} \quad (51)$$

Tuning of λ : a SURE-based approach

Recall:

$$\|\mathbf{y} - \widehat{\mathbf{y}}\|^2 + 2\text{tr}(S\Sigma) \text{ is an unbiased estimator of the risk } \mathbb{E}\left[\|\mathbf{y}^* - \widehat{\mathbf{y}}\|_2^2\right] \text{ with } \widehat{\mathbf{y}} = S\mathbf{y} \quad (51)$$

In our case:

$$S = \frac{G^T G}{M} \left(\frac{G^T G}{M} + \frac{\gamma\sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \quad (52)$$

$$\Sigma = \sigma_\nu^2 \frac{G^T G}{M^2}. \quad (53)$$

Tuning of λ : a SURE-based approach

Recall:

$$\|\mathbf{y} - \widehat{\mathbf{y}}\|^2 + 2\text{tr}(S\Sigma) \text{ is an unbiased estimator of the risk } \mathbb{E} \left[\|\mathbf{y}^* - \widehat{\mathbf{y}}\|_2^2 \right] \text{ with } \widehat{\mathbf{y}} = S\mathbf{y} \quad (51)$$

In our case:

$$S = \frac{G^T G}{M} \left(\frac{G^T G}{M} + \frac{\gamma \sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \quad (52)$$

$$\Sigma = \sigma_\nu^2 \frac{G^T G}{M^2}. \quad (53)$$

Important consideration: the original process was $\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu}$, but this SURE approach considers only $G\mathbf{a} + \boldsymbol{\nu}$!

Tuning of λ : a SURE-based approach

Recall:

$$\|\mathbf{y} - \widehat{\mathbf{y}}\|^2 + 2\text{tr}(S\Sigma) \text{ is an unbiased estimator of the risk } \mathbb{E} \left[\|\mathbf{y}^* - \widehat{\mathbf{y}}\|_2^2 \right] \text{ with } \widehat{\mathbf{y}} = S\mathbf{y} \quad (51)$$

In our case:

$$S = \frac{G^T G}{M} \left(\frac{G^T G}{M} + \frac{\gamma \sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \quad (52)$$

$$\Sigma = \sigma_\nu^2 \frac{G^T G}{M^2}. \quad (53)$$

Important consideration: the original process was $\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu}$, but this SURE approach considers only $G\mathbf{a} + \boldsymbol{\nu}$! However, for large M , $G^T Z\mathbf{b}$ vanishes \implies SURE score above is an asymptotically unbiased estimator of the actual risk

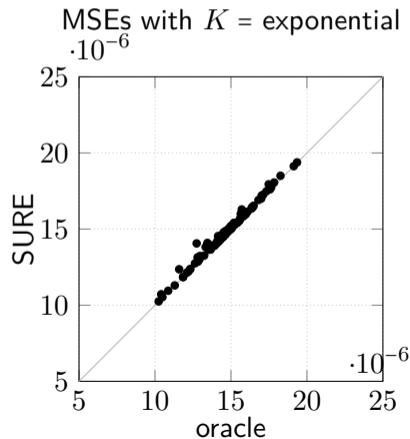
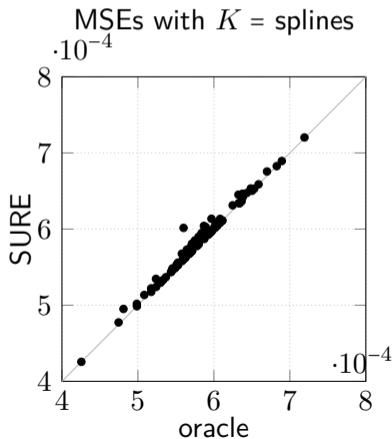
Does this work? Analysis on synthetic data

Does this work? Analysis on synthetic data

- $M = 10000$, 1000 Monte-Carlo runs, $K =$ splines or exponential
- $\Lambda = 50$ potential λ s, log-spaced in $[10^{-3}, 10^3]$
- \exists “oracle” that knows f and thus what is the best λ

Does this work? Analysis on synthetic data

- $M = 10000$, 1000 Monte-Carlo runs, $K =$ splines or exponential
- $\Lambda = 50$ potential λ s, log-spaced in $[10^{-3}, 10^3]$
- \exists “oracle” that knows f and thus what is the best λ



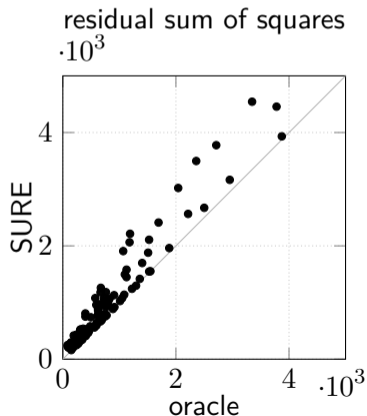
Does this work? Analysis on field data

Does this work? Analysis on field data

- “Colorado rain” UCI dataset
- $K(x, x') = \exp(-10 \|x - x'\|_2^2)$ fixed, Λ as before
- 1000 Monte-Carlo runs, each with 2 random months of data as training set and 1 random month as test

Does this work? Analysis on field data

- “Colorado rain” UCI dataset
- $K(x, x') = \exp(-10 \|x - x'\|_2^2)$ fixed, Λ as before
- 1000 Monte-Carlo runs, each with 2 random months of data as training set and 1 random month as test



Some brief concluding remarks

Some brief concluding remarks

- SURE approaches can be valid alternatives to other model selection strategies
- seem to be suitable for distributed average-consensus based settings

Some brief concluding remarks

- SURE approaches can be valid alternatives to other model selection strategies
- seem to be suitable for distributed average-consensus based settings

What now?

- time-varying estimation
- generalizations for other distributed estimation settings & big data
(not discussed in this presentation)

Purposes of this seminar

- 1 discuss about a useful tool
- 2 connect with you