

Exponentially Weighted Aggregate with the Laplace Prior

Joint work with Edwin Grappin and Quentin Paris



3rd UCL Workshop on the Theory of Big Data
Mon 26-Wed 28 June 2017

Arnak S. Dalalyan

ENSAE ParisTech / CREST

Estimation of sparse vectors and matrices

Problem formulation

Regression model and prediction

- **Regression model:** we observe n feature-label pairs $(\mathbf{X}_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}$ such that

$$Y_i = \mathbf{X}_i^\top \boldsymbol{\beta}^* + \xi_i \Leftrightarrow \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\xi}$$

with $\boldsymbol{\xi} = [\xi_1, \dots, \xi_n] \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

- **Regression vector:** The p -dimensional vector $\boldsymbol{\beta}^*$ is unknown.
- **Deterministic design:** The vectors \mathbf{X}_i are deterministic.

Problem formulation

Regression model and prediction

- **Regression model:** we observe n feature-label pairs $(\mathbf{X}_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}$ such that

$$Y_i = \mathbf{X}_i^\top \boldsymbol{\beta}^* + \xi_i \Leftrightarrow \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\xi}$$

with $\boldsymbol{\xi} = [\xi_1, \dots, \xi_n] \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

- **Regression vector:** The p -dimensional vector $\boldsymbol{\beta}^*$ is unknown.
- **Deterministic design:** The vectors \mathbf{X}_i are deterministic.
- **Prediction loss:** The quality of an estimator $\hat{\boldsymbol{\beta}}_n$ is measured by

$$\ell_n(\hat{\boldsymbol{\beta}}_n, \boldsymbol{\beta}^*) = \frac{1}{n} \|\mathbf{X}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*)\|_2^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_n - \mathbf{X}_i^\top \boldsymbol{\beta}^*)^2.$$

- **Curse of dimensionality:** OLS behaves poorly if $n \ll p$.

Problem formulation

Trace regression

- **Trace Regression:** we observe n pairs $(\mathbf{X}_i, Y_i) \in \mathcal{M}_{m_1, m_2} \times \mathbb{R}$ such that

$$Y_i = \langle \mathbf{X}_i, \mathbf{B}^* \rangle + \xi_i$$

with $\boldsymbol{\xi} = [\xi_1, \dots, \xi_n] \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ and $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{A}^\top \mathbf{B})$.

- The $m_1 \times m_2$ -dimensional matrix \mathbf{B}^* is unknown.
- **Deterministic design:** The matrices \mathbf{X}_i are deterministic.

Problem formulation

Trace regression

- **Trace Regression:** we observe n pairs $(\mathbf{X}_i, Y_i) \in \mathcal{M}_{m_1, m_2} \times \mathbb{R}$ such that

$$Y_i = \langle \mathbf{X}_i, \mathbf{B}^* \rangle + \xi_i$$

with $\boldsymbol{\xi} = [\xi_1, \dots, \xi_n] \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ and $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{A}^\top \mathbf{B})$.

- The $m_1 \times m_2$ -dimensional matrix \mathbf{B}^* is unknown.
- **Deterministic design:** The matrices \mathbf{X}_i are deterministic.
- **Prediction loss:** The quality of an estimator $\widehat{\mathbf{B}}_n$ is measured by

$$\ell_n(\widehat{\mathbf{B}}_n, \mathbf{B}^*) = \frac{1}{n} \sum_{i=1}^n (\langle \mathbf{X}_i, \widehat{\mathbf{B}}_n \rangle - \langle \mathbf{X}_i, \mathbf{B}^* \rangle)^2.$$

Sparsity, approximate sparsity and mis-specification

- **s -sparsity:** When p is large, to make consistent estimation of β^* feasible, we assume that

$$\|\beta^*\|_0 = \#\{j : \beta_j^* \neq 0\} =: s \ll n \wedge p.$$

Sparsity, approximate sparsity and mis-specification

- **s -sparsity:** When p is large, to make consistent estimation of β^* feasible, we assume that

$$\|\beta^*\|_0 = \#\{j : \beta_j^* \neq 0\} =: s \ll n \wedge p.$$

- **δ -approximate s -sparsity:** for some $\bar{\beta} \in \mathbb{R}^p$,

$$\|\bar{\beta}\|_0 \leq s \quad \text{and} \quad \|\bar{\beta} - \beta^*\|_1 \leq \delta.$$

Sparsity, approximate sparsity and mis-specification

- **s -sparsity:** When p is large, to make consistent estimation of β^* feasible, we assume that

$$\|\beta^*\|_0 = \#\{j : \beta_j^* \neq 0\} =: s \ll n \wedge p.$$

- **δ -approximate s -sparsity:** for some $\bar{\beta} \in \mathbb{R}^p$,

$$\|\bar{\beta}\|_0 \leq s \quad \text{and} \quad \|\bar{\beta} - \beta^*\|_1 \leq \delta.$$

- **Mis-specification:** β^* is neither sparse nor approximately sparse, but there is an approximately sparse vector $\bar{\beta}$ with a small prediction error $\ell_n(\bar{\beta}, \beta^*)$.

This $\bar{\beta}$ is the (δ -approximately s -sparse) oracle.
Can we mimic the oracle ?

Low-rank, approximate low-rank and mis-specification

- The role of sparsity in matrix estimation problems is most often played by the low rank assumption (sparsity of the spectrum).
- **Matrix sparsity:** \mathbf{B}^* is r -low-rank, if

$$\text{rank}(\mathbf{B}^*) \leq r \ll m_1 \wedge m_2.$$

Low-rank, approximate low-rank and mis-specification

- The role of sparsity in matrix estimation problems is most often played by the low rank assumption (sparsity of the spectrum).
- **Matrix sparsity:** \mathbf{B}^* is r -low-rank, if

$$\text{rank}(\mathbf{B}^*) \leq r \ll m_1 \wedge m_2.$$

- δ -approximate r -low rank: for some $\bar{\mathbf{B}}$,

$$\text{rank}(\bar{\mathbf{B}}) \leq r$$

and

$$\|\bar{\mathbf{B}} - \mathbf{B}^*\|_1 \leq \delta.$$

Low-rank, approximate low-rank and mis-specification

- The role of sparsity in matrix estimation problems is most often played by the low rank assumption (sparsity of the spectrum).
- **Matrix sparsity:** \mathbf{B}^* is r -low-rank, if

$$\text{rank}(\mathbf{B}^*) \leq r \ll m_1 \wedge m_2.$$

- δ -approximate r -low rank: for some $\bar{\mathbf{B}}$,

$$\text{rank}(\bar{\mathbf{B}}) \leq r$$

and

$$\|\bar{\mathbf{B}} - \mathbf{B}^*\|_1 \leq \delta.$$

nuclear norm = sum of singular values

Low-rank, approximate low-rank and mis-specification

- The role of sparsity in matrix estimation problems is most often played by the low rank assumption (sparsity of the spectrum).
- **Matrix sparsity:** \mathbf{B}^* is r -low-rank, if

$$\text{rank}(\mathbf{B}^*) \leq r \ll m_1 \wedge m_2.$$

- δ -approximate r -low rank: for some $\bar{\mathbf{B}}$,

$$\text{rank}(\bar{\mathbf{B}}) \leq r$$

and

$$\|\bar{\mathbf{B}} - \mathbf{B}^*\|_1 \leq \delta.$$

nuclear norm = sum of singular values

- **Mis-specification:** there is an approximately low rank matrix $\bar{\mathbf{B}}$ with a small prediction error $\ell_n(\bar{\mathbf{B}}, \mathbf{B}^*)$.

Lasso, Bayesian Lasso and EWA

Estimating sparse vectors by the lasso

- **Regression model, prediction loss:** $n \times p$ matrix \mathbf{X} and $\mathbf{Y} \in \mathbb{R}^p$ such that

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\xi},$$

$$\ell_n(\hat{\boldsymbol{\beta}}_n, \boldsymbol{\beta}^*) = 1/n \|\mathbf{X}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*)\|_2^2.$$

Estimating sparse vectors by the lasso

- **Regression model, prediction loss:** $n \times p$ matrix \mathbf{X} and $\mathbf{Y} \in \mathbb{R}^p$ such that

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\xi},$$

$$\ell_n(\widehat{\boldsymbol{\beta}}_n, \boldsymbol{\beta}^*) = 1/n \|\mathbf{X}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*)\|_2^2.$$

- **Weak sparsity:** Good prediction can be achieved by a nearly sparse vector $\widehat{\boldsymbol{\beta}}$.

Estimating sparse vectors by the lasso

- **Regression model, prediction loss:** $n \times p$ matrix \mathbf{X} and $\mathbf{Y} \in \mathbb{R}^p$ such that

$$\mathbf{Y} = \mathbf{X}\beta^* + \xi,$$

$$\ell_n(\hat{\beta}_n, \beta^*) = 1/n \|\mathbf{X}(\hat{\beta}_n - \beta^*)\|_2^2.$$

- **Weak sparsity:** Good prediction can be achieved by a nearly sparse vector $\hat{\beta}$.
- **Lasso [Tibshirani 96]:** ℓ_1 penalized least squares

$$\hat{\beta}^L \in \arg \min_{\beta \in \mathbb{R}^p} \{1/2n \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1.\}$$

- Fast computation by LARS or convex programming.

Prediction accuracy of the lasso

- **Lasso [Tibshirani 1996]:** ℓ_1 penalized least squares

$$\hat{\beta}^L \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ 1/2n \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

- **Oracle inequality** [Koltchinskii, Lounici and Tsybakov 2011; D., Hebiri and Lederer 2017]: if $\lambda \geq \sqrt{(2/n) \log(p/\alpha)}$ then with probability $\geq 1 - \alpha$

$$l_n(\hat{\beta}^L, \beta^*) \leq l_n(\bar{\beta}, \beta^*) + \min_J \left\{ 4\lambda \|\bar{\beta}_{J^c}\|_1 + (3/\kappa)\lambda^2 |J| \right\}.$$

Prediction accuracy of the lasso

- **Lasso [Tibshirani 1996]:** ℓ_1 penalized least squares

$$\hat{\beta}^L \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

- **Oracle inequality** [Koltchinskii, Lounici and Tsybakov 2011; D., Hebiri and Lederer 2017]: if $\lambda \geq \sqrt{(2/n) \log(p/\alpha)}$ then with probability $\geq 1 - \alpha$

$$\ell_n(\hat{\beta}^L, \beta^*) \leq \ell_n(\bar{\beta}, \beta^*) + \min_J \left\{ 4\lambda \|\bar{\beta}_{J^c}\|_1 + (3/\kappa)\lambda^2 |J| \right\}.$$

- **Slow rate:** For $J = \emptyset$ the rate is $\sqrt{1/n}$ without κ .

Prediction accuracy of the lasso

- **Lasso [Tibshirani 1996]:** ℓ_1 penalized least squares

$$\hat{\beta}^L \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

- **Oracle inequality** [Koltchinskii, Lounici and Tsybakov 2011; D., Hebiri and Lederer 2017]: if $\lambda \geq \sqrt{(2/n) \log(p/\alpha)}$ then with probability $\geq 1 - \alpha$

$$\ell_n(\hat{\beta}^L, \beta^*) \leq \ell_n(\bar{\beta}, \beta^*) + \min_J \left\{ 4\lambda \|\bar{\beta}_{J^c}\|_1 + (3/\kappa)\lambda^2 |J| \right\}.$$

- **Slow rate:** For $J = \emptyset$ the rate is $\sqrt{1/n}$ without κ .
- **Fast rate:** For getting the fast rate $(s/n) \log(p)$, we do need the restricted eigenvalue (RE) condition $\kappa > 0$ [Bickel, Ritov, Tsybakov (2009), D., Hebiri and Lederer 2017].

Estimating sparse vectors by the Bayesian lasso

- **Regression with Gaussian noise:**

$$Y = \mathbf{X}\beta^* + \xi, \quad \xi \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_n).$$

Estimating sparse vectors by the Bayesian lasso

- **Regression with Gaussian noise:**

$$Y = \mathbf{X}\beta^* + \xi, \xi \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_n).$$

- **Laplace prior:** sparsity inducing prior

$$\pi_0(\beta) \propto \exp\left(- (n\lambda/\sigma^2) \|\beta\|_1\right).$$

Estimating sparse vectors by the Bayesian lasso

- **Regression with Gaussian noise:**

$$Y = \mathbf{X}\beta^* + \xi, \xi \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_n).$$

- **Laplace prior:** sparsity inducing prior

$$\pi_0(\beta) \propto \exp\left(- (n\lambda/\sigma^2) \|\beta\|_1\right).$$

- **Posterior density** is then

$$\hat{\pi}_n(\beta) \propto \exp\left\{- (1/2\sigma^2) \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 - (n\lambda/\sigma^2) \|\beta\|_1\right\}.$$

- **Bayesian lasso** [Park and Casella 2008]: posterior mean

$$\hat{\beta}^{\text{BL}} = \int \beta \hat{\pi}_n(d\beta).$$

- **Remark:** the MAP coincides with the lasso.

Properties of the Bayesian lasso

● Pros

- Computable by the Gibbs sampling or Langevin-type algorithm.
- Uncertainty quantification ?
- Less sensitive to changes in data than the lasso.

Properties of the Bayesian lasso

● Pros

- Computable by the Gibbs sampling or Langevin-type algorithm.
- Uncertainty quantification ?
- Less sensitive to changes in data than the lasso.

● Cons

- Computational guarantees weaker than for the lasso.
- No matter the λ , $\hat{\beta}^{\text{BL}}$ does not converge at the optimal rate [Castillo, Schmidt-Hieber and van der Vaart 2015], $\forall C > 0$

$$\hat{\pi}_n \left(\left\{ \beta : \ell_n(\beta, \beta^*) \leq \frac{Cs \log p}{n} \right\} \right) \xrightarrow[n \rightarrow \infty]{} 0, \quad \text{a.s.}$$

even when the matrix \mathbf{X} is diagonal.

Properties of the Bayesian lasso

● Pros

- Computable by the Gibbs sampling or Langevin-type algorithm.
- Uncertainty quantification ?
- Less sensitive to changes in data than the lasso.

● Cons

- Computational guarantees weaker than for the lasso.
- No matter the λ , $\hat{\beta}^{\text{BL}}$ does not converge at the optimal rate [Castillo, Schmidt-Hieber and van der Vaart 2015], $\forall C > 0$

$$\hat{\pi}_n \left(\left\{ \beta : \ell_n(\beta, \beta^*) \leq \frac{Cs \log p}{n} \right\} \right) \xrightarrow[n \rightarrow \infty]{} 0, \quad \text{a.s.}$$

even when the matrix \mathbf{X} is diagonal.

One solution is to change the prior (spike-and-slab) but then the computation becomes NP-hard.

EWA and PAC-Bayes approach

- **Pseudo-posterior:** for some prior π_0 , define

$$\hat{\pi}_{n,\tau} \in \arg \min_{p(\cdot)} \left\{ \int_{\mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 p(d\boldsymbol{\beta}) + \tau \mathcal{K}(p \parallel \pi_0) \right\},$$

where $\tau > 0$ is a tuning parameter referred to as temperature.

EWA and PAC-Bayes approach

- **Pseudo-posterior:** for some prior π_0 , define

$$\hat{\pi}_{n,\tau} \in \arg \min_{p(\cdot)} \left\{ \int_{\mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 p(d\boldsymbol{\beta}) + \tau \mathcal{K}(p \parallel \pi_0) \right\},$$

where $\tau > 0$ is a tuning parameter referred to as temperature.

- **Explicit form:**

$$\hat{\pi}_{n,\tau}(\boldsymbol{\beta}) \propto \exp \left\{ -\frac{1}{2n\tau} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \log \pi_0(\boldsymbol{\beta}) \right\}.$$

EWA and PAC-Bayes approach

- **Pseudo-posterior:** for some prior π_0 , define

$$\hat{\pi}_{n,\tau} \in \arg \min_{p(\cdot)} \left\{ \int_{\mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 p(d\boldsymbol{\beta}) + \tau \mathcal{K}(p \parallel \pi_0) \right\},$$

where $\tau > 0$ is a tuning parameter referred to as temperature.

- **Explicit form:**

$$\hat{\pi}_{n,\tau}(\boldsymbol{\beta}) \propto \exp \left\{ -\frac{1}{2n\tau} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \log \pi_0(\boldsymbol{\beta}) \right\}.$$

- **Exponentially weighted aggregate (EWA):**

$$\hat{\boldsymbol{\beta}}_{\tau}^{\text{EWA}} = \int_{\mathbb{R}^p} \boldsymbol{\beta} \hat{\pi}_{n,\tau}(d\boldsymbol{\beta}).$$

- If $\tau = \sigma^2/n$, then the EWA coincides with the Bayes estimator.

EWA and the Laplace prior

- One can use a scaled Laplace prior in the EWA.
- The resulting (suitably re-parameterized) pseudo-posterior is

$$\hat{\pi}_{n,\tau}(\boldsymbol{\beta}) \propto \exp \left\{ -\frac{1}{2n\tau} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 - \frac{\lambda}{\tau} \|\boldsymbol{\beta}\|_1 \right\}.$$

- This estimator interpolates the lasso and the Bayesian lasso:
 - For $\tau = 0$, the EWA coincides with the lasso.
 - For $\tau = \sigma^2/n$, the EWA coincides with the Bayesian lasso.

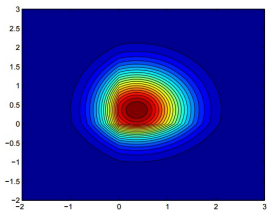
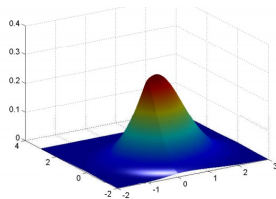
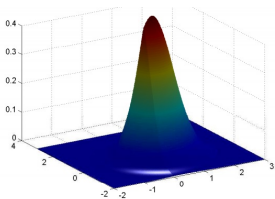
EWA and the Laplace prior

- One can use a scaled Laplace prior in the EWA.
- The resulting (suitably re-parameterized) pseudo-posterior is

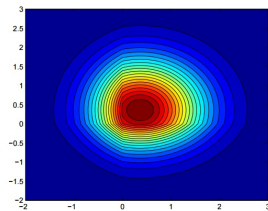
$$\hat{\pi}_{n,\tau}(\boldsymbol{\beta}) \propto \exp \left\{ -\frac{1}{2n\tau} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 - \frac{\lambda}{\tau} \|\boldsymbol{\beta}\|_1 \right\}.$$

- This estimator interpolates the lasso and the Bayesian lasso:
 - For $\tau = 0$, the EWA coincides with the lasso.
 - For $\tau = \sigma^2/n$, the EWA coincides with the Bayesian lasso.
- **Questions:**
 - Is the EWA suitable for prediction under the sparsity scenario? If it is, what is the range of temperature τ providing good prediction accuracy?
 - How do the statistical properties of the EWA compare with those of the lasso?

EWA and the Laplace prior



$\tau = 0.5$



$\tau = 0.8$

Main results

Oracle inequality for the EWA

Theorem 1

Assume that $\xi \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and that $\max_{j \in [p]} 1/n \|\mathbf{x}^j\|_2^2 \leq 1$. Suppose, in addition, that $\lambda \geq 2\sigma \sqrt{(2/n) \log(p/\alpha)}$, for some $\alpha \in (0, 1)$. Then, with probability at least $1 - \alpha$,

$$l_n(\widehat{\beta}_\tau^{\text{EWA}}, \beta^*) \leq \min_{J \subset [p]} \left\{ l_n(\bar{\beta}, \beta^*) + 4\lambda \|\bar{\beta}_{J^c}\|_1 + \frac{9\lambda^2 |J|}{4\kappa} \right\} + 2p\tau. \quad (1)$$

Oracle inequality for the EWA

Theorem 1

Assume that $\xi \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and that $\max_{j \in [p]} 1/n \|\mathbf{x}^j\|_2^2 \leq 1$. Suppose, in addition, that $\lambda \geq 2\sigma \sqrt{(2/n) \log(p/\alpha)}$, for some $\alpha \in (0, 1)$. Then, with probability at least $1 - \alpha$,

$$\ell_n(\widehat{\beta}_\tau^{\text{EWA}}, \beta^*) \leq \min_{J \subset [p]} \left\{ \ell_n(\bar{\beta}, \beta^*) + 4\lambda \|\bar{\beta}_{J^c}\|_1 + \frac{9\lambda^2 |J|}{4\kappa} \right\} + 2p\tau. \quad (1)$$

Remarks

- 1 This is a sharp OI, since the leading constant is 1.
- 2 This risk bound extends the risk bounds available for the lasso [D., Hebiri and Lederer 2017] to the EWA.
- 3 The last term in the above risk bound, $2p\tau$, reflects the influence of the temperature τ ; if $\tau = \sigma^2/(pn)$ this term is negligible.

Theorem 2

Assume that $\xi \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and that $\max_{j \in [p]} 1/n \|\mathbf{x}^j\|_2^2 \leq 1$. Assume that $\lambda \geq 2\sigma \sqrt{(2/n) \log(p/\alpha)}$. Then, with probability at least $1 - \alpha$ on data, the pseudo-posterior $\hat{\pi}_{n,\tau}$ with the Laplace prior puts at least a weight $1 - 2e^{-\sqrt{p}/16}$ on the set

$$\left\{ \beta : \ell_n(\beta, \beta^*) \leq \min_{J \subset [p]} \left(\ell_n(\bar{\beta}, \beta^*) + 4\lambda \|\bar{\beta}_{J^c}\|_1 + \frac{9\lambda^2 |J|}{2\kappa} \right) + 8p\tau \right\}.$$

Pseudo-posterior concentration

Theorem 2

Assume that $\xi \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and that $\max_{j \in [p]} 1/n \|\mathbf{x}^j\|_2^2 \leq 1$. Assume that $\lambda \geq 2\sigma \sqrt{(2/n) \log(p/\alpha)}$. Then, with probability at least $1 - \alpha$ on data, the pseudo-posterior $\hat{\pi}_{n,\tau}$ with the Laplace prior puts at least a weight $1 - 2e^{-\sqrt{p}/16}$ on the set

$$\left\{ \beta : \ell_n(\beta, \beta^*) \leq \min_{J \subset [p]} \left(\ell_n(\bar{\beta}, \beta^*) + 4\lambda \|\bar{\beta}_{J^c}\|_1 + \frac{9\lambda^2 |J|}{2\kappa} \right) + 8p\tau \right\}.$$

Remarks

- 1 A random sample from $\hat{\pi}_{n,\tau}$ is almost as good as the pseudo-posterior mean.
- 2 It might be easier to sample from $\hat{\pi}_{n,\tau}$ rather than to compute the pseudo-posterior mean.
- 3 The method is sparsity-adaptive.

Main idea of the proof of Theorem 1

- Let $V_n(\beta)$ be the penalized log-likelihood.
- The lasso estimator and the EWA are

$$\hat{\beta}^L \in \arg \min_{\beta} V_n(\beta), \quad \hat{\beta}_{\tau}^{\text{EWA}} = \frac{1}{C} \int_{\mathbb{R}^p} \beta \exp\{-V_n(\beta)/\tau\} d\beta.$$

- **Known:** the proof of the OI for the lasso builds on

$$\nabla V_n(\hat{\beta}^L) = \mathbf{0} \quad [\text{rigorously } \mathbf{0} \in \partial V_n(\hat{\beta}^L)].$$

- **New:** the analogue of the above condition is:

$$\int_{\mathbb{R}^p} \nabla V_n(\beta) \hat{\pi}_n(\beta) d\beta = C \int_{\mathbb{R}^p} \nabla [\hat{\pi}_n(\beta)] d\beta = \mathbf{0}.$$

- Additionally, for every $j \in [p]$, $\int_{\mathbb{R}^p} \nabla [\beta_j \hat{\pi}_n(\beta)] d\beta = \mathbf{0}$.

Main ingredient of the proof of Theorem 2

Magic theorem [Bobkov and Madiman 2011]

Let $\hat{\pi}_n(\mathbf{u}) \propto \exp(-V_n(\mathbf{u})/\tau)$ be a log-concave probability density^a and let β be a random vector drawn from $\hat{\pi}_n$. Then, for any $t > 0$, the inequality

$$V_n(\beta) \leq \int_{\mathbb{R}^p} V_n(\mathbf{u}) \hat{\pi}_n(\mathbf{u}) d\mathbf{u} + \sqrt{p\tau t}$$

holds with probability at least $1 - 2e^{-t/16}$.

^aThis means that V_n is a convex function.

Lasso versus EWA

- Both are computable in polynomial time, but the lasso computation is faster.
- Theoretical guarantees for the lasso and the EWA are equally good, provided that the temperature τ is small.
- The EWA is less sensitive to the variations in the data.
 - The SURE for the lasso [Tibshirani and Taylor, 2012] is given by

$$\widehat{R}^L(\lambda) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}^L(\lambda)\|_2^2 - \sigma^2 + \frac{2\sigma^2}{n} \text{rank}(\mathbf{X}_{\mathcal{A}(\lambda)}),$$

where $\mathcal{A}(\lambda) = \{j \in [p] : \beta_j^L(\lambda) \neq 0\}$.

This function is even not continuous.

- The SURE for the EWA is

$$\widehat{R}^{\text{EWA}}(\lambda) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau}^{\text{EWA}}(\lambda)\|_2^2 - \sigma^2 + \frac{2\sigma^2}{n^2\tau} \mathbf{Var}_{\widehat{\pi}_n(\lambda)}[\mathbf{X}\boldsymbol{\beta}].$$

OI for the EWA in the matrix trace regression

Recall that the ℓ_1 -norm $\|\beta\|_1$ is replaced by the nuclear norm $\|\mathbf{B}\|_1$.

Let us define $v_{\mathbf{X}} = \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right\|^{1/2} \vee \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^\top \mathbf{X}_i \right\|^{1/2}$.

Theorem 3

Assume that $\lambda \geq 2\sigma v_{\mathbf{X}} \{2/n \log((m_1 + m_2)/\alpha)\}^{1/2}$, for some $\alpha \in (0, 1)$. Then, with probability at least $1 - \alpha$, the matrix $\hat{\mathbf{B}}^{\text{EWA}}$ satisfies

$$\ell_n(\hat{\mathbf{B}}^{\text{EWA}}, \mathbf{B}^*) \leq \min_J \left\{ \ell_n(\bar{\mathbf{B}}, \mathbf{B}^*) + 4\lambda \|\mathcal{P}_{J^c}(\bar{\mathbf{B}})\|_1 + \frac{9\lambda^2 |J|}{4\kappa} \right\} + 2m_1 m_2 \tau,$$

where the min is over all subsets $J \subset [\text{rank}(\bar{\mathbf{B}})]$.

We have a similar result on pseudo-posterior concentration.

Summary and outlook

- We have obtained nonasymptotic risk bounds for the EWA with Laplace prior, both in vector- and matrix-regression models.
- If the temperature $\tau \leq \sigma^2/(np)$, then the EWA is as good as the lasso.
- The pseudo-posterior has nice concentration properties.
- This can be extended to random design, but more technical and additional terms appear.
- Current/future work:
 - Computational complexity of drawing from $\hat{\pi}_{n,\tau}$.
 - Assessing the sensitivity of the EWA as compared to the lasso.

