

# Matrix superharmonic priors for Bayes estimation under matrix quadratic loss

Takeru Matsuda (RIKEN Center for Brain Science)

William E. Strawderman (Rutgers University)

# Abstract

## Stein (1974)

When  $X \sim N_n(\mu, I_n)$  ( $n \geq 3$ ), Bayes estimator with a superharmonic prior  $\pi(\mu)$  is minimax under quadratic loss:

$$\Delta\pi := \sum_{a=1}^n \frac{\partial^2 \pi}{\partial \mu_a^2} \leq 0 \quad \Rightarrow \quad E\|\hat{\mu}^\pi(x) - \mu\|^2 \leq n$$

## This study (M. and Strawderman, *Biometrika* 2021+)

When  $X \sim N_{n,p}(M, I_n, I_p)$  ( $n \geq p + 2$ ), Bayes estimator with a matrix superharmonic prior is minimax under matrix quadratic loss:

$$\begin{aligned} \tilde{\Delta}\pi &:= \left( \sum_{a=1}^n \frac{\partial^2 \pi}{\partial M_{ai} \partial M_{aj}} \right)_{ij} \leq O \\ \Rightarrow \quad E(\hat{M}^\pi(X) - M)^\top (\hat{M}^\pi(X) - M) &\leq nI_p \end{aligned}$$

# Stein's paradox

$$X \sim N_n(\mu, I_n)$$

- estimate  $\mu$  based on  $X$  under quadratic loss  $\|\hat{\mu} - \mu\|^2$
- Maximum likelihood estimator  $\hat{\mu}_{\text{MLE}}(x) = x$  is minimax.

## Theorem (Stein, 1956)

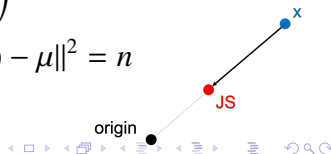
When  $n \geq 3$ ,  $\hat{\mu}_{\text{MLE}}(x) = x$  is inadmissible.

- **Shrinkage estimators** dominate  $\hat{\mu}_{\text{MLE}}$ .
- e.g. James–Stein estimator (James and Stein, 1961)

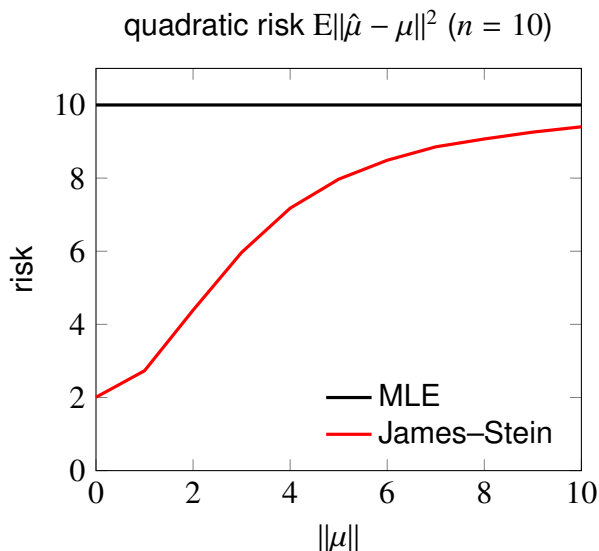
$$\hat{\mu}_{\text{JS}}(x) = \left(1 - \frac{n-2}{\|x\|^2}\right)x$$

$$E\|\hat{\mu}_{\text{JS}}(x) - \mu\|^2 \leq E\|\hat{\mu}_{\text{MLE}}(x) - \mu\|^2 = n$$

- JS shrinks  $x$  toward the origin.



## Risk comparison



- JS attains large risk reduction when  $\mu$  is close to the origin

## superharmonic prior $\Rightarrow$ minimax

- Bayes estimator of  $\mu$  with prior  $\pi(\mu)$  (posterior mean)

$$\hat{\mu}^\pi(x) = \mathbb{E}_\pi[\mu | x] = \int \mu \pi(\mu | x) d\mu$$

- **superharmonic** prior

$$\Delta\pi(\mu) = \sum_{a=1}^n \frac{\partial^2}{\partial \mu_a^2} \pi(\mu) \leq 0$$

### Theorem (Stein, 1974)

The Bayes estimator with a superharmonic prior is minimax.

- e.g. Stein's prior ( $n \geq 3$ )

$$\pi_S(\mu) = \|\mu\|^{2-n}$$

- Bayes estimator with  $\pi_S$  shrinks toward the origin like JS.

# Shrinkage estimation of normal mean matrix

$$X \sim N_{n,p}(M, I_n, I_p) \quad (X_{ai} \sim N(M_{ai}, 1))$$

- estimate  $M$  based on  $X$  under Frobenius loss

$$\|\hat{M} - M\|_F^2 = \sum_{a=1}^n \sum_{i=1}^p (\hat{M}_{ai} - M_{ai})^2$$

- Efron–Morris estimator (= James–Stein estimator when  $p = 1$ )

$$\hat{M}_{EM}(X) = X \left( I_p - (n - p - 1)(X^T X)^{-1} \right)$$

## Theorem (Efron and Morris, 1972)

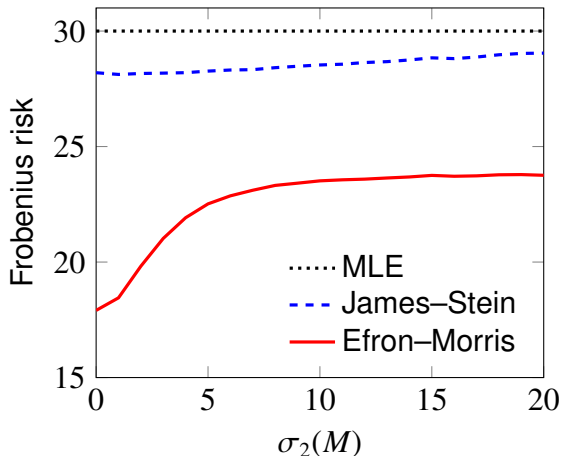
When  $n \geq p + 2$ ,  $\hat{M}_{EM}$  is minimax and dominates  $\hat{M}_{MLE}(X) = X$ .

- Stein (1974):  $\hat{M}_{EM}$  **shrinks singular values** separately.

$$\sigma_i(\hat{M}_{EM}) = \left( 1 - \frac{n - p - 1}{\sigma_i(X)^2} \right) \sigma_i(X)$$

## Risk function (rank 2)

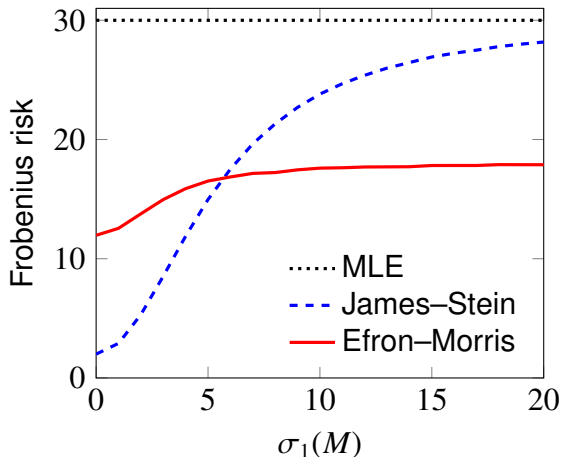
- $n = 10, p = 3, \sigma_1(M) = 20, \sigma_3(M) = 0$



- $\hat{M}_{EM}$  works well when  $\sigma_2(M)$  is small, **even if  $\sigma_1(M)$  is large.**
  - ▶  $\hat{M}_{JS}$  works well if  $\|M\|_F^2 = \sigma_1(M)^2 + \sigma_2(M)^2 + \sigma_3(M)^2$  is small.

## Risk function (rank 1)

- $n = 10, p = 3, \sigma_2(M) = \sigma_3(M) = 0$



- $\hat{M}_{EM}$  has constant risk reduction even if  $\sigma_1(M)$  is large.
- Therefore,  $\hat{M}_{EM}$  works well when  $M$  has **low rank**.



## Singular value shrinkage prior (M. and Komaki, 2015)

$$\pi_{\text{SVS}}(M) = \det(M^T M)^{-(n-p-1)/2} = \prod_{i=1}^p \sigma_i(M)^{-(n-p-1)}$$

- puts more weight on matrices with smaller singular values  
→ **shrinks singular values separately**
- When  $p = 1$ ,  $\pi_{\text{SVS}}$  coincides with Stein's prior  $\pi_S(\mu) = \|\mu\|^{2-n}$ .

### Theorem (M. and Komaki, *Biometrika* 2015)

When  $n \geq p + 2$ ,  $\pi_{\text{SVS}}$  is superharmonic:

$$\Delta \pi_{\text{SVS}} = \sum_{a=1}^n \sum_{i=1}^p \frac{\partial^2 \pi_{\text{SVS}}}{\partial M_{ai}^2} \leq 0.$$

- Bayes estimator with  $\pi_{\text{SVS}}$  is minimax under Frobenius loss.
  - ▶ similar behavior to EM
  - ▶ works well when  $M$  has (approximately) low rank

## Summary (so far)

vector	matrix
James–Stein estimator $\hat{\mu}_{\text{JS}} = \left(1 - \frac{n-2}{\ x\ ^2}\right)x$	Efron–Morris estimator $\hat{M}_{\text{EM}} = X \left( I_p - (n-p-1)(X^T X)^{-1} \right)$
Stein's prior $\pi_{\text{S}}(\mu) = \ \mu\ ^{-(n-2)}$	singular value shrinkage prior $\pi_{\text{SVS}}(M) = \det(M^T M)^{-(n-p-1)/2}$

- note: JS and EM are not (generalized) Bayes estimators.

# Estimation under matrix quadratic loss

$$X \sim \mathcal{N}_{n,p}(M, I_n, I_p) \quad (X_{ai} \sim \mathcal{N}(M_{ai}, 1))$$

- estimate  $M$  based on  $X$  under **matrix quadratic loss**

$$L(M, \hat{M}) = (\hat{M} - M)^\top (\hat{M} - M) \in \mathbb{R}^{p \times p}$$

- risk function

$$R(M, \hat{M}) = \mathbb{E}_M[L(M, \hat{M}(X))] \in \mathbb{R}^{p \times p}$$

- We compare  $R(M, \hat{M})$  in the **Löwner order**  $\leq$ 
  - ▶  $A \leq B \Leftrightarrow B - A$  is positive semidefinite
  - ▶ If  $R(M, \hat{M}_1) \leq R(M, \hat{M}_2)$ , then  $\mathbb{E}_M \|\hat{M}_1 c - M c\|^2 \leq \mathbb{E}_M \|\hat{M}_2 c - M c\|^2$  for every  $c$
- cf. multivariate linear regression

## Unbiased risk estimate & minimaxity of EM

- matrix divergence

$$(\widetilde{\text{div}} g(X))_{ij} = \sum_{a=1}^n \frac{\partial}{\partial X_{ai}} g_{aj}(X)$$

### Theorem

The matrix quadratic risk of  $\hat{M} = X + g(X)$  is given by

$$R(M, \hat{M}) = nI_p + E_M[\widetilde{\text{div}} g(X) + (\widetilde{\text{div}} g(X))^T + g(X)^T g(X)]$$

### Theorem

When  $n - p - 1 > 0$ , the Efron–Morris estimator is minimax under the matrix quadratic loss:

$$R(M, \hat{M}_{\text{EM}}) = nI_p - (n - p - 1)^2 E_M[(X^T X)^{-1}] \leq nI_p$$

## Matrix superharmonicity

- “sphere” with center  $X \in \mathbb{R}^{n \times p}$  and “radius”  $\rho \in \mathbb{R}^p$

$$S_{X,\rho} = \{X + e\rho^\top \mid e \in \mathbb{R}^n, \|e\| = 1\}$$

- average value of  $f$  on  $S_{X,\rho}$

$$L(f : X, \rho) = \frac{1}{\Omega_n} \int_{S_{0,1}} f(X + e\rho^\top) ds(e)$$

### Definition

An extended real-valued function  $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R} \cup \{\infty\}$  is matrix superharmonic if

- 1  $f$  is lower semicontinuous
- 2  $f \not\equiv \infty$
- 3  $L(f : X, \rho) \leq f(X)$  for every  $X \in \mathbb{R}^{n \times p}$  and  $\rho \in \mathbb{R}^p$

# Matrix superharmonic $\Rightarrow$ superharmonic

## Lemma

If a function  $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R} \cup \{\infty\}$  is matrix superharmonic, then  $f \circ \text{vec}^{-1}$  is superharmonic.

- Proof: For every  $X \in \mathbb{R}^{n \times p}$  and  $r > 0$ ,

$$L(f \circ \text{vec}^{-1} : \text{vec}(X), r) = \frac{1}{\Omega_p r^{p-1}} \int_{S_{0,r}} L(f : X, \rho) ds(\rho) \leq f(X)$$

- The converse does not hold when  $p \geq 2$ .
  - ▶ e.g.  $f(X) = \|X\|_F^{2-np}$

## Characterization by matrix Laplacian

- Matrix superharmonicity is characterized by a matrix version of the Laplacian.

### Definition

For a  $C^2$  function  $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$ , its matrix Laplacian  $\tilde{\Delta}f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{p \times p}$  is defined as

$$(\tilde{\Delta}f(X))_{ij} = \sum_{a=1}^n \frac{\partial^2}{\partial X_{ai} \partial X_{aj}} f(X)$$

### Theorem

A  $C^2$  function  $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$  is matrix superharmonic if and only if its matrix Laplacian is negative semidefinite  $\tilde{\Delta}f(X) \leq O$  for every  $X$ .

- Proof: Green's theorem

## matrix superharmonic prior $\Rightarrow$ minimax

$$\hat{M}^\pi(X) = E_\pi[M | X] = X + \tilde{\nabla} \log m_\pi(X)$$

### Theorem

If  $\sqrt{m_\pi(X)}$  is matrix superharmonic, then  $\hat{M}^\pi$  is minimax under the matrix quadratic loss.

- Proof: by using the unbiased estimate of risk,

$$R(M, \hat{M}^\pi) = nI_p + 4E_M \left[ \frac{\tilde{\Delta} \sqrt{m_\pi(X)}}{\sqrt{m_\pi(X)}} \right]$$

### Theorem

If  $\pi(M)$  is matrix superharmonic, then  $\sqrt{m_\pi(X)}$  is also matrix superharmonic and  $\hat{M}^\pi$  is minimax under the matrix quadratic loss.

- When  $p = 1$ , it reduces to the classical result by Stein (1974).



# A class of matrix superharmonic priors

- improper matrix t-prior

$$\pi_{\alpha,\beta}(M) = \det(M^T M + \beta I_p)^{-(\alpha+n+p-1)/2}$$

## Theorem

If  $-n - p + 1 \leq \alpha \leq -2p$  and  $\beta \geq 0$ , then  $\pi_{\alpha,\beta}(M)$  is matrix superharmonic and the generalized Bayes estimator with respect to  $\pi_{\alpha,\beta}(M)$  is minimax under the matrix quadratic loss.

- When  $p = 1$ , it reduces to the result by Faith (1993) on (improper) multivariate t-priors.

## Matrix superharmonicity of $\pi_{\text{SVS}}$

- When  $\alpha = -2p$  and  $\beta = 0$ , the prior  $\pi_{\alpha,\beta}(M)$  coincides with the singular value shrinkage prior

$$\pi_{\text{SVS}}(M) = \det(M^T M)^{-(n-p-1)/2}$$

### Corollary

When  $n - p - 1 > 0$ ,  $\pi_{\text{SVS}}(M)$  is matrix superharmonic and the generalized Bayes estimator with respect to  $\pi_{\text{SVS}}$  is minimax under the matrix quadratic loss.

- The matrix superharmonicity of  $\pi_{\text{SVS}}$  is strongly concentrated on the space of low rank matrices.

### Corollary

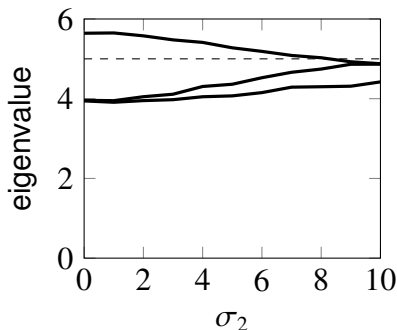
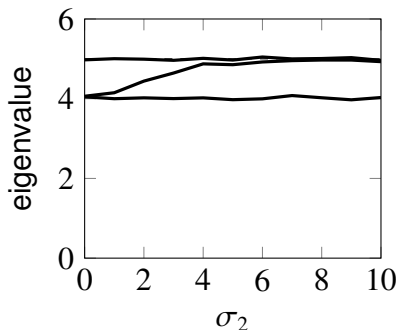
If  $M$  has full-rank, then  $\tilde{\Delta}\pi_{\text{SVS}}(M) = O$ .

## Simulation setting

- We denote the  $i$ -th singular value of  $M$  by  $\sigma_i$ .
  - ▶  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$
- We focus on the eigenvalues  $\lambda_1 \geq \dots \geq \lambda_p$  of the matrix quadratic risk  $R(M, \hat{M})$ .
  - ▶ Since  $R(M, \hat{M}) = nI_p$  for MLE  $\hat{M} = X$ , an estimator is minimax if and only if  $\lambda_1 \leq n$  for every  $M$ .
- Bayes estimator with  $\pi_{\text{SVS}}(M) = \det(M^\top M)^{-(n-p-1)/2}$
- Bayes estimator with Stein's prior  $\pi_{\text{S}}(M) = \|M\|_{\text{F}}^{2-np}$
- Efron–Morris estimator  $\hat{M}_{\text{EM}} = X(I - (n - p - 1)(X^\top X)^{-1})$ 
  - ▶ almost the same risk with Bayes estimators with  $\pi_{\text{SVS}}$
- James–Stein estimator  $\hat{M}_{\text{JS}} = (1 - (np - 2)/\|X\|_{\text{F}}^2)X$ 
  - ▶ almost the same risk with Bayes estimators with  $\pi_{\text{S}}$

## Simulation results (Figure 1)

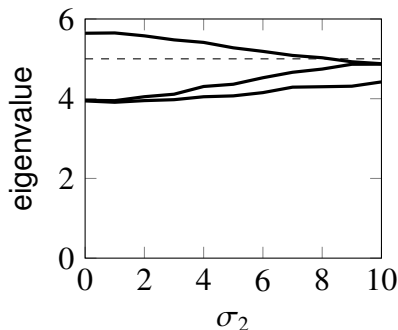
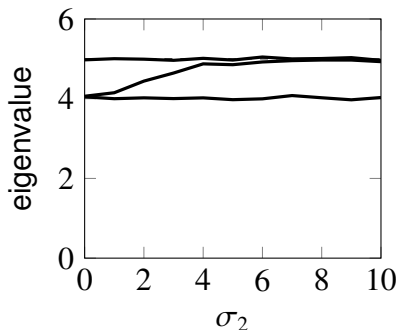
- eigenvalues of Bayes estimators ( $n = 5, p = 3, \sigma_1 = 10, \sigma_3 = 0$ )
- left:  $\pi_{SVS}$ , right:  $\pi_S$



- For  $\pi_{SVS}$ , all eigenvalues do not exceed  $n = 5$ , which indicates the minimaxity.
  - ▶  $\lambda_1$  and  $\lambda_3$  are almost constant with values  $\lambda_1 \approx 5$  and  $\lambda_3 \approx 4$ .
  - ▶  $\lambda_2$  increases from 4 to 5 with  $\sigma_2$ .
  - ▶ These are understood from the fact that  $\pi_{SVS}$  shrinks each singular value separately.

## Simulation results (Figure 1)

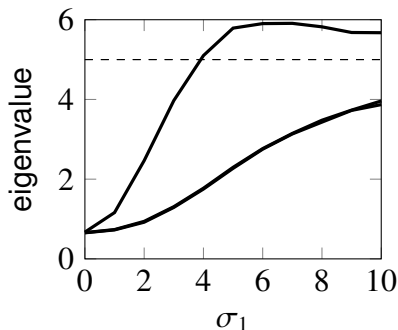
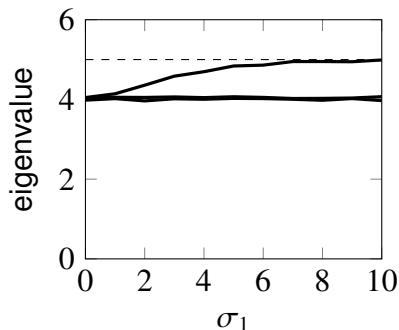
- eigenvalues of Bayes estimators ( $n = 5, p = 3, \sigma_1 = 10, \sigma_3 = 0$ )
- left:  $\pi_{SVS}$ , right:  $\pi_S$



- For  $\pi_S$ ,  $\lambda_1 \geq n = 5$  when  $\sigma_2 \leq 8 \rightarrow$  not minimax.
  - ▶ However, since this estimator is minimax under the Frobenius loss,  $\lambda_1 + \lambda_2 + \lambda_3 \leq np = 15$ .
  - ▶ cf. James–Stein estimator is not minimax componentwise, even though it is minimax under the quadratic loss for the whole vector (Lehmann and Casella, 2006).

## Simulation results (Figure 2)

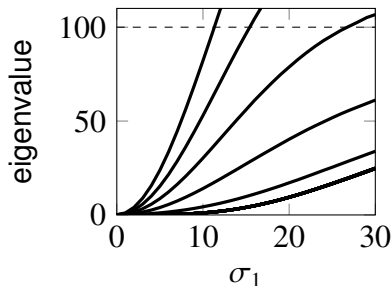
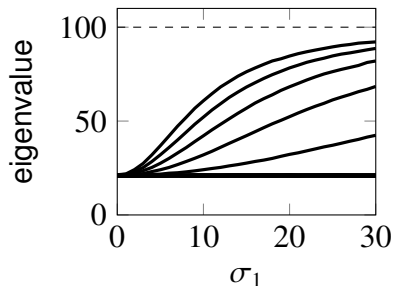
- eigenvalues ( $n = 5$ ,  $p = 3$ ,  $\sigma_2 = \sigma_3 = 0$ )
- left:  $\pi_{\text{SVS}}$ , right:  $\pi_{\text{S}}$



- For  $\pi_{\text{SVS}}$ , both  $\lambda_2$  and  $\lambda_3$  are almost constant around 4  
→  $\pi_{\text{SVS}}$  works particularly well when  $M$  has low rank

## Simulation results (Figure 3)

- eigenvalues ( $n = 100$ ,  $p = 20$ ,  $\sigma_i = (6 - i)/5 \cdot \sigma_1$  ( $i = 2, \dots, 5$ ),  $\sigma_6 = \dots = \sigma_{20} = 0$ )
- left:  $\hat{M}_{EM}$ . right:  $\hat{M}_{JS}$



- The advantage of  $\hat{M}_{EM}$  to the low-rank setting is more pronounced in higher dimensions.
  - ▶  $\lambda_6 \approx \dots \approx \lambda_{20} \approx 20$

# Summary

$$X \sim N_{n,p}(M, I_n, I_p)$$

- The Bayes estimator with a **matrix superharmonic** prior is minimax under **matrix quadratic loss**:

$$\tilde{\Delta}\pi := \left( \sum_{a=1}^n \frac{\partial^2 \pi}{\partial M_{ai} \partial M_{aj}} \right)_{ij} \leq O$$

$$\Rightarrow E(\hat{M}^\pi(X) - M)^\top (\hat{M}^\pi(X) - M) \leq nI_p$$

- The matrix t-prior

$$\pi_{\alpha,\beta}(M) = \det(M^\top M + \beta I_p)^{-(\alpha+n+p-1)/2}$$

is matrix superharmonic when  $-n - p + 1 \leq \alpha \leq -2p$  and  $\beta \geq 0$ .

- Matrix superharmonic priors work well for **low-rank** matrices.
- paper: M. and Strawderman, *Biometrika* 2021+