

# High Dimensional Predictive Inference

Ed George  
University of Pennsylvania  
(joint work with Larry Brown, Feng Liang, and Xinyi Xu)

*Conference on Predictive Inference and Its Applications*  
*Iowa State University*  
*May 7, 2018*



# I. The Hunt for Shrinkage Estimators Begins

# I. The Hunt for Shrinkage Estimators Begins

- ▶ Canonical Problem: Observe  $X \mid \mu \sim N_p(\mu, I)$  and estimate  $\mu$  by  $\hat{\mu}$  under

$$R_Q(\mu, \hat{\mu}) = E_\mu \|\hat{\mu}(X) - \mu\|^2$$

# I. The Hunt for Shrinkage Estimators Begins

- ▶ Canonical Problem: Observe  $X \mid \mu \sim N_p(\mu, I)$  and estimate  $\mu$  by  $\hat{\mu}$  under

$$R_Q(\mu, \hat{\mu}) = E_{\mu} \|\hat{\mu}(X) - \mu\|^2$$

- ▶  $\hat{\mu}_{MLE}(X) = X$  is MLE, best invariant and minimax with constant risk  $R_Q(\mu, \hat{\mu}_{MLE}) \equiv p$ .

# I. The Hunt for Shrinkage Estimators Begins

- ▶ Canonical Problem: Observe  $X \mid \mu \sim N_p(\mu, I)$  and estimate  $\mu$  by  $\hat{\mu}$  under

$$R_Q(\mu, \hat{\mu}) = E_\mu \|\hat{\mu}(X) - \mu\|^2$$

- ▶  $\hat{\mu}_{MLE}(X) = X$  is MLE, best invariant and minimax with constant risk  $R_Q(\mu, \hat{\mu}_{MLE}) \equiv p$ .
- ▶ A Shocking Discovery:  $\hat{\mu}_{MLE}$  is inadmissible when  $p \geq 3$ . (Stein 1956)

## I. The Hunt for Shrinkage Estimators Begins

- ▶ Canonical Problem: Observe  $X \mid \mu \sim N_p(\mu, I)$  and estimate  $\mu$  by  $\hat{\mu}$  under

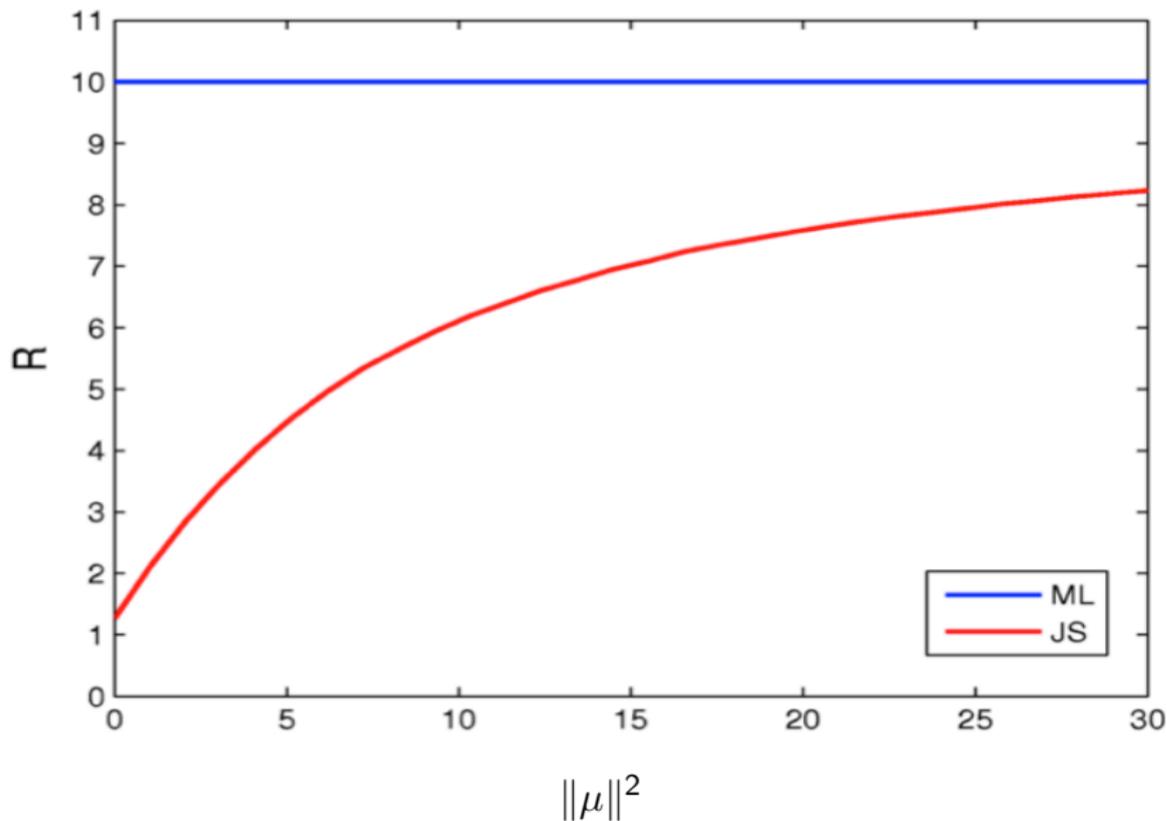
$$R_Q(\mu, \hat{\mu}) = E_\mu \|\hat{\mu}(X) - \mu\|^2$$

- ▶  $\hat{\mu}_{MLE}(X) = X$  is MLE, best invariant and minimax with constant risk  $R_Q(\mu, \hat{\mu}_{MLE}) \equiv p$ .
- ▶ A Shocking Discovery:  $\hat{\mu}_{MLE}$  is inadmissible when  $p \geq 3$ . (Stein 1956)
- ▶ An Explicit Better Estimator Appears: The James-Stein estimator

$$\hat{\mu}_{JS} = \left(1 - \frac{p-2}{\|X\|^2}\right) X$$


(James and Stein 1961)

- The risk of  $\hat{\mu}_{MLE}$  and the risk of  $\hat{\mu}_{JS}$  various values of  $\mu$



- ▶ Stein (1962) suggests an empirical Bayes motivation for  $\hat{\mu}_{JS}$ .  
The focus of the hunt turns to Bayes.

- ▶ Stein (1962) suggests an empirical Bayes motivation for  $\hat{\mu}_{JS}$ . The focus of the hunt turns to Bayes.
- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_Q(\mu, \hat{\mu})$  is

$$\hat{\mu}_\pi(X) = E_\pi(\mu | X)$$

- ▶ Stein (1962) suggests an empirical Bayes motivation for  $\hat{\mu}_{JS}$ . The focus of the hunt turns to Bayes.
- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_Q(\mu, \hat{\mu})$  is

$$\hat{\mu}_\pi(X) = E_\pi(\mu | X)$$

- ▶ Remark: The (formal) Bayes rule under  $\pi_U(\mu) \equiv 1$  is

$$\hat{\mu}_U(X) \equiv \hat{\mu}_{MLE}(X) = X$$

- ▶  $\hat{\mu}_H(X)$ , the Bayes rule under the Harmonic prior

$$\pi_H(\mu) = \|\mu\|^{-(p-2)},$$

dominates  $\hat{\mu}_U$  when  $p \geq 3$ . (Stein 1974)

- ▶  $\hat{\mu}_H(X)$ , the Bayes rule under the Harmonic prior

$$\pi_H(\mu) = \|\mu\|^{-(p-2)},$$

dominates  $\hat{\mu}_U$  when  $p \geq 3$ . (Stein 1974)

- ▶  $\hat{\mu}_a(X)$ , the Bayes rule under  $\pi_a(\mu)$  where

$$\mu \mid s \sim N_p(0, sI), \quad s \sim (1+s)^{a-2}$$

dominates  $\hat{\mu}_U$  and is proper Bayes when  $p = 5$  and  $a \in [.5, 1)$  or when  $p \geq 6$  and  $a \in [0, 1)$ . (Strawderman 1971)

- ▶  $\hat{\mu}_H(X)$ , the Bayes rule under the Harmonic prior

$$\pi_H(\mu) = \|\mu\|^{-(p-2)},$$

dominates  $\hat{\mu}_U$  when  $p \geq 3$ . (Stein 1974)

- ▶  $\hat{\mu}_a(X)$ , the Bayes rule under  $\pi_a(\mu)$  where

$$\mu \mid s \sim N_p(0, sI), \quad s \sim (1+s)^{a-2}$$

dominates  $\hat{\mu}_U$  and is proper Bayes when  $p = 5$  and  $a \in [.5, 1)$  or when  $p \geq 6$  and  $a \in [0, 1)$ . (Strawderman 1971)

- ▶ A Unifying Phenomenon: These domination results can be attributed to properties of the marginal distribution of  $X$  under  $\pi_H$  and  $\pi_a$ . 

- ▶ The Bayes rule under  $\pi(\mu)$  can be expressed as

$$\hat{\mu}_\pi(X) = E_\pi(\mu | X) = X + \nabla \log m_\pi(X)$$

where

$$m_\pi(X) \propto \int e^{-(X-\mu)^2/2} \pi(\mu) d\mu$$

is the marginal of  $X$  under  $\pi(\mu)$ . ( $\nabla = (\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_p})'$ )  
(Brown 1971)

- ▶ The Bayes rule under  $\pi(\mu)$  can be expressed as

$$\hat{\mu}_\pi(X) = E_\pi(\mu | X) = X + \nabla \log m_\pi(X)$$

where

$$m_\pi(X) \propto \int e^{-(X-\mu)^2/2} \pi(\mu) d\mu$$

is the marginal of  $X$  under  $\pi(\mu)$ . ( $\nabla = (\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_p})'$ )  
(Brown 1971)

- ▶ The risk improvement of  $\hat{\mu}_\pi(X)$  over  $\hat{\mu}_U(X)$  can be expressed as

$$\begin{aligned} R_Q(\mu, \hat{\mu}_U) - R_Q(\mu, \hat{\mu}_\pi) &= E_\mu \left[ \|\nabla \log m_\pi(X)\|^2 - 2 \frac{\nabla^2 m_\pi(X)}{m_\pi(X)} \right] \\ &= E_\mu \left[ -4 \nabla^2 \sqrt{m_\pi(X)} / \sqrt{m_\pi(X)} \right] \end{aligned}$$

( $\nabla^2 = \sum_i \frac{\partial^2}{\partial x_i^2}$ ) (Stein 1974, 1981)

- ▶ That  $\hat{\mu}_H(X)$  dominates  $\hat{\mu}_U$  when  $p \geq 3$ , follows from the fact that the marginal  $m_\pi(X)$  under  $\pi_H$  is superharmonic, i.e.

$$\nabla^2 m_\pi(X) \leq 0$$

- ▶ That  $\hat{\mu}_H(X)$  dominates  $\hat{\mu}_U$  when  $p \geq 3$ , follows from the fact that the marginal  $m_\pi(X)$  under  $\pi_H$  is superharmonic, i.e.

$$\nabla^2 m_\pi(X) \leq 0$$

- ▶ That  $\hat{\mu}_a(X)$  dominates  $\hat{\mu}_U$  when  $p \geq 5$  (and conditions on  $a$ ), follows from the fact that the sqrt of the marginal under  $\pi_a$  is superharmonic, i.e.

$$\nabla^2 \sqrt{m_\pi(X)} \leq 0$$

(Fourdrinier, Strawderman and Wells 1998)



## II. The Prediction Problem

## II. The Prediction Problem

- ▶ Observe  $X \mid \mu \sim N_p(\mu, v_x I)$  and predict  $Y \mid \mu \sim N_p(\mu, v_y I)$ 
  - ▶ Given  $\mu$ ,  $Y$  is independent of  $X$
  - ▶  $v_x$  and  $v_y$  are known (for now)

## II. The Prediction Problem

- ▶ Observe  $X | \mu \sim N_p(\mu, v_x I)$  and predict  $Y | \mu \sim N_p(\mu, v_y I)$ 
  - ▶ Given  $\mu$ ,  $Y$  is independent of  $X$
  - ▶  $v_x$  and  $v_y$  are known (for now)
- ▶ The Problem: To estimate  $p(y | \mu)$  by  $q(y | x)$ .

## II. The Prediction Problem

- ▶ Observe  $X | \mu \sim N_p(\mu, v_x I)$  and predict  $Y | \mu \sim N_p(\mu, v_y I)$ 
  - ▶ Given  $\mu$ ,  $Y$  is independent of  $X$
  - ▶  $v_x$  and  $v_y$  are known (for now)
- ▶ The Problem: To estimate  $p(y | \mu)$  by  $q(y | x)$ .
- ▶ Measure closeness by Kullback-Leibler loss,

$$L(\mu, q(y | x)) = \int p(y | \mu) \log \frac{p(y | \mu)}{q(y | x)} dy$$

## II. The Prediction Problem

- ▶ Observe  $X | \mu \sim N_p(\mu, v_x I)$  and predict  $Y | \mu \sim N_p(\mu, v_y I)$ 
  - ▶ Given  $\mu$ ,  $Y$  is independent of  $X$
  - ▶  $v_x$  and  $v_y$  are known (for now)
- ▶ The Problem: To estimate  $p(y | \mu)$  by  $q(y | x)$ .
- ▶ Measure closeness by Kullback-Leibler loss,

$$L(\mu, q(y | x)) = \int p(y | \mu) \log \frac{p(y | \mu)}{q(y | x)} dy$$

- ▶ Risk function

$$R_{KL}(\mu, q) = \int L(\mu, q(y | x)) p(x | \mu) dx = E_{\mu}[L(\mu, q(y | X))] \quad \text{⏪ ⏩ 🔍 ↺}$$

# Bayes Rules for the Prediction Problem

## Bayes Rules for the Prediction Problem

- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_{KL}(\mu, q)$  is

$$p_{\pi}(y | x) = \int p(y | \mu)\pi(\mu | x)d\mu = E_{\pi}[p(y | \mu)|X]$$

## Bayes Rules for the Prediction Problem

- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_{KL}(\mu, q)$  is

$$p_{\pi}(y | x) = \int p(y | \mu)\pi(\mu | x)d\mu = E_{\pi}[p(y | \mu)|X]$$

- ▶ Let  $p_U(y | x)$  denote the (formal) Bayes rule under  $\pi_U(\mu) \equiv 1$ .

## Bayes Rules for the Prediction Problem

- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_{KL}(\mu, q)$  is

$$p_{\pi}(y | x) = \int p(y | \mu)\pi(\mu | x)d\mu = E_{\pi}[p(y | \mu)|X]$$

- ▶ Let  $p_U(y | x)$  denote the (formal) Bayes rule under  $\pi_U(\mu) \equiv 1$ .
- ▶  $p_U(y | x)$  dominates  $p(y | \hat{\mu} = x)$ , the naive “plug-in” predictive distribution. (Aitchison 1975)

## Bayes Rules for the Prediction Problem

- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_{KL}(\mu, q)$  is

$$p_{\pi}(y | x) = \int p(y | \mu)\pi(\mu | x)d\mu = E_{\pi}[p(y | \mu)|X]$$

- ▶ Let  $p_U(y | x)$  denote the (formal) Bayes rule under  $\pi_U(\mu) \equiv 1$ .
- ▶  $p_U(y | x)$  dominates  $p(y | \hat{\mu} = x)$ , the naive “plug-in” predictive distribution. (Aitchison 1975)
- ▶  $p_U(y | x)$  is best invariant and minimax with constant risk. (Murray 1977, Ng 1980, Barron and Liang 2003)

## Bayes Rules for the Prediction Problem

- ▶ For a prior  $\pi(\mu)$ , the Bayes rule under  $R_{KL}(\mu, q)$  is

$$p_{\pi}(y | x) = \int p(y | \mu)\pi(\mu | x)d\mu = E_{\pi}[p(y | \mu)|X]$$

- ▶ Let  $p_U(y | x)$  denote the (formal) Bayes rule under  $\pi_U(\mu) \equiv 1$ .
- ▶  $p_U(y | x)$  dominates  $p(y | \hat{\mu} = x)$ , the naive “plug-in” predictive distribution. (Aitchison 1975)
- ▶  $p_U(y | x)$  is best invariant and minimax with constant risk. (Murray 1977, Ng 1980, Barron and Liang 2003)
- ▶ Shocking Fact:  $p_U(y | x)$  is inadmissible when  $p \geq 3$ .

- ▶  $p_H(y | x)$ , the Bayes rule under the Harmonic prior

$$\pi_H(\mu) = \|\mu\|^{-(p-2)},$$

dominates  $p_U(y | x)$  when  $p \geq 3$ . (Komaki 2001)

- ▶  $p_H(y | x)$ , the Bayes rule under the Harmonic prior

$$\pi_H(\mu) = \|\mu\|^{-(p-2)},$$

dominates  $p_U(y | x)$  when  $p \geq 3$ . (Komaki 2001)

- ▶  $p_a(y | x)$ , the Bayes rule under  $\pi_a(\mu)$  where

$$\mu | s \sim N_p(0, s v_0 I), \quad s \sim (1 + s)^{a-2},$$

dominates  $p_U(y | x)$  and is proper Bayes when  $v_x \leq v_0$  and when  $p = 5$  and  $a \in [.5, 1)$  or when  $p \geq 6$  and  $a \in [0, 1)$ . (Liang 2002)

- ▶  $p_H(y | x)$ , the Bayes rule under the Harmonic prior

$$\pi_H(\mu) = \|\mu\|^{-(p-2)},$$

dominates  $p_U(y | x)$  when  $p \geq 3$ . (Komaki 2001)

- ▶  $p_a(y | x)$ , the Bayes rule under  $\pi_a(\mu)$  where

$$\mu | s \sim N_p(0, s v_0 I), \quad s \sim (1 + s)^{a-2},$$

dominates  $p_U(y | x)$  and is proper Bayes when  $v_x \leq v_0$  and when  $p = 5$  and  $a \in [.5, 1)$  or when  $p \geq 6$  and  $a \in [0, 1)$ . (Liang 2002)

- ▶ A Key Question: Are these domination results attributable to the properties of  $m_\pi$ ?



## A Key Representation for $p_{\pi}(y | x)$

## A Key Representation for $p_\pi(y | x)$

- ▶ Let  $m_\pi(x; v_x)$  denote the marginal of  $X | \mu \sim N_p(\mu, v_x I)$  under  $\pi(\mu)$ .

## A Key Representation for $p_\pi(y | x)$

- ▶ Let  $m_\pi(x; v_x)$  denote the marginal of  $X | \mu \sim N_p(\mu, v_x I)$  under  $\pi(\mu)$ .
- ▶ **Lemma:** The Bayes rule  $p_\pi(y | x)$  can be expressed as

$$p_\pi(y | x) = \frac{m_\pi(w; v_w)}{m_\pi(x; v_x)} p_U(y | x)$$

where

$$W = \frac{v_y X + v_x Y}{v_x + v_y} \sim N_p(\mu, v_w I)$$

## A Key Representation for $p_\pi(y | x)$

- ▶ Let  $m_\pi(x; v_x)$  denote the marginal of  $X | \mu \sim N_p(\mu, v_x I)$  under  $\pi(\mu)$ .
- ▶ **Lemma:** The Bayes rule  $p_\pi(y | x)$  can be expressed as

$$p_\pi(y | x) = \frac{m_\pi(w; v_w)}{m_\pi(x; v_x)} p_U(y | x)$$

where

$$W = \frac{v_y X + v_x Y}{v_x + v_y} \sim N_p(\mu, v_w I)$$

- ▶ Using this, the risk improvement can be expressed as

$$\begin{aligned} R_{KL}(\mu, p_U) - R_{KL}(\mu, p_\pi) &= \int \int p_{v_x}(x | \mu) p_{v_y}(y | \mu) \log \frac{p_\pi(y | x)}{p_U(y | x)} dx dy \\ &= E_{\mu, v_w} \log m_\pi(W; v_w) - E_{\mu, v_x} \log m_\pi(X; v_x) \end{aligned}$$

# An Analogue of Stein's Unbiased Estimate of Risk

## An Analogue of Stein's Unbiased Estimate of Risk

► **Theorem:**

$$\begin{aligned}\frac{\partial}{\partial \mathbf{v}} E_{\mu, \mathbf{v}} \log m_{\pi}(Z; \mathbf{v}) &= E_{\mu, \mathbf{v}} \left[ \frac{\nabla^2 m_{\pi}(Z; \mathbf{v})}{m_{\pi}(Z; \mathbf{v})} - \frac{1}{2} \|\nabla \log m_{\pi}(Z; \mathbf{v})\|^2 \right] \\ &= E_{\mu, \mathbf{v}} \left[ 2\nabla^2 \sqrt{m_{\pi}(Z; \mathbf{v})} / \sqrt{m_{\pi}(Z; \mathbf{v})} \right]\end{aligned}$$

## An Analogue of Stein's Unbiased Estimate of Risk

► **Theorem:**

$$\begin{aligned}\frac{\partial}{\partial \mathbf{v}} E_{\mu, \mathbf{v}} \log m_{\pi}(Z; \mathbf{v}) &= E_{\mu, \mathbf{v}} \left[ \frac{\nabla^2 m_{\pi}(Z; \mathbf{v})}{m_{\pi}(Z; \mathbf{v})} - \frac{1}{2} \|\nabla \log m_{\pi}(Z; \mathbf{v})\|^2 \right] \\ &= E_{\mu, \mathbf{v}} \left[ 2\nabla^2 \sqrt{m_{\pi}(Z; \mathbf{v})} / \sqrt{m_{\pi}(Z; \mathbf{v})} \right]\end{aligned}$$

► Proof relies on using the heat equation

$$\frac{\partial}{\partial \mathbf{v}} m_{\pi}(z; \mathbf{v}) = \frac{1}{2} \nabla^2 m_{\pi}(z; \mathbf{v}),$$



Brown's representation and Stein's Lemma.

# General Conditions for Minimax Prediction

## General Conditions for Minimax Prediction

- ▶ Let  $m_\pi(z; \nu)$  be the marginal distribution of  $Z \mid \mu \sim N_p(\mu, \nu I)$  under  $\pi(\mu)$ .

## General Conditions for Minimax Prediction

- ▶ Let  $m_\pi(z; \nu)$  be the marginal distribution of  $Z \mid \mu \sim N_p(\mu, \nu I)$  under  $\pi(\mu)$ .
- ▶ **Theorem:** If  $m_\pi(z; \nu)$  is finite for all  $z$ , then  $p_\pi(y \mid x)$  will be minimax if either of the following hold:
  1.  $m_\pi(z; \nu)$  is superharmonic
  2.  $\sqrt{m_\pi(z; \nu)}$  is superharmonic

## General Conditions for Minimax Prediction

- ▶ Let  $m_\pi(z; \nu)$  be the marginal distribution of  $Z \mid \mu \sim N_p(\mu, \nu I)$  under  $\pi(\mu)$ .
- ▶ **Theorem:** If  $m_\pi(z; \nu)$  is finite for all  $z$ , then  $p_\pi(y \mid x)$  will be minimax if either of the following hold:
  1.  $m_\pi(z; \nu)$  is superharmonic
  2.  $\sqrt{m_\pi(z; \nu)}$  is superharmonic
- ▶ **Corollary:** If  $m_\pi(z; \nu)$  is finite for all  $z$ , then  $p_\pi(y \mid x)$  will be minimax if  $\pi(\mu)$  is superharmonic.

## General Conditions for Minimax Prediction

- ▶ Let  $m_\pi(z; \nu)$  be the marginal distribution of  $Z \mid \mu \sim N_p(\mu, \nu I)$  under  $\pi(\mu)$ .
- ▶ **Theorem:** If  $m_\pi(z; \nu)$  is finite for all  $z$ , then  $p_\pi(y \mid x)$  will be minimax if either of the following hold:
  1.  $m_\pi(z; \nu)$  is superharmonic
  2.  $\sqrt{m_\pi(z; \nu)}$  is superharmonic
- ▶ **Corollary:** If  $m_\pi(z; \nu)$  is finite for all  $z$ , then  $p_\pi(y \mid x)$  will be minimax if  $\pi(\mu)$  is superharmonic.
- ▶  $p_\pi(y \mid x)$  will dominate  $p_U(y \mid x)$  in the above results if the superharmonicity is strict on some interval.

# Consequences of the General Minimax Conditions

## Consequences of the General Minimax Conditions

- ▶ Because  $\pi_H$  is superharmonic, it is immediate that  $p_H(y | x)$  dominates  $p_U(y | x)$  and is minimax.

## Consequences of the General Minimax Conditions

- ▶ Because  $\pi_H$  is superharmonic, it is immediate that  $p_H(y | x)$  dominates  $p_U(y | x)$  and is minimax.
- ▶ Because  $\sqrt{m_a}$  is superharmonic (under suitable conditions on  $a$ ), it is immediate that  $p_a(y | x)$  dominates  $p_U(y | x)$  and is minimax.

## Consequences of the General Minimax Conditions

- ▶ Because  $\pi_H$  is superharmonic, it is immediate that  $p_H(y | x)$  dominates  $p_U(y | x)$  and is minimax.
- ▶ Because  $\sqrt{m_a}$  is superharmonic (under suitable conditions on  $a$ ), it is immediate that  $p_a(y | x)$  dominates  $p_U(y | x)$  and is minimax.
- ▶ It also follows that any of the improper superharmonic t-priors of Faith (1978) or any of the proper generalized t-priors of Fourdrinier, Strawderman and Wells (1998) yield Bayes rules that dominate  $p_U(y | x)$  and are minimax.

### III. Predictive “Shrinkage”

### III. Predictive “Shrinkage”

- ▶ Our Lemma representation

$$p_H(y | x) = \frac{m_H(w; v_w)}{m_H(x; v_x)} p_U(y | x)$$

shows how  $p_H(y | x)$  “shrinks  $p_U(y | x)$  towards 0” by an adaptive multiplicative factor.

### III. Predictive “Shrinkage”

- ▶ Our Lemma representation

$$p_H(y | x) = \frac{m_H(w; v_w)}{m_H(x; v_x)} p_U(y | x)$$

shows how  $p_H(y | x)$  “shrinks  $p_U(y | x)$  towards 0” by an adaptive multiplicative factor.

- ▶ Note the analogies with the Bayes rule  $\hat{\mu}_\pi(X) = E_\pi(\mu | X)$  whose coordinates are

$$\left( 1 + \frac{(\nabla \log m_\pi(X))_i}{X_i} \right) X_i$$

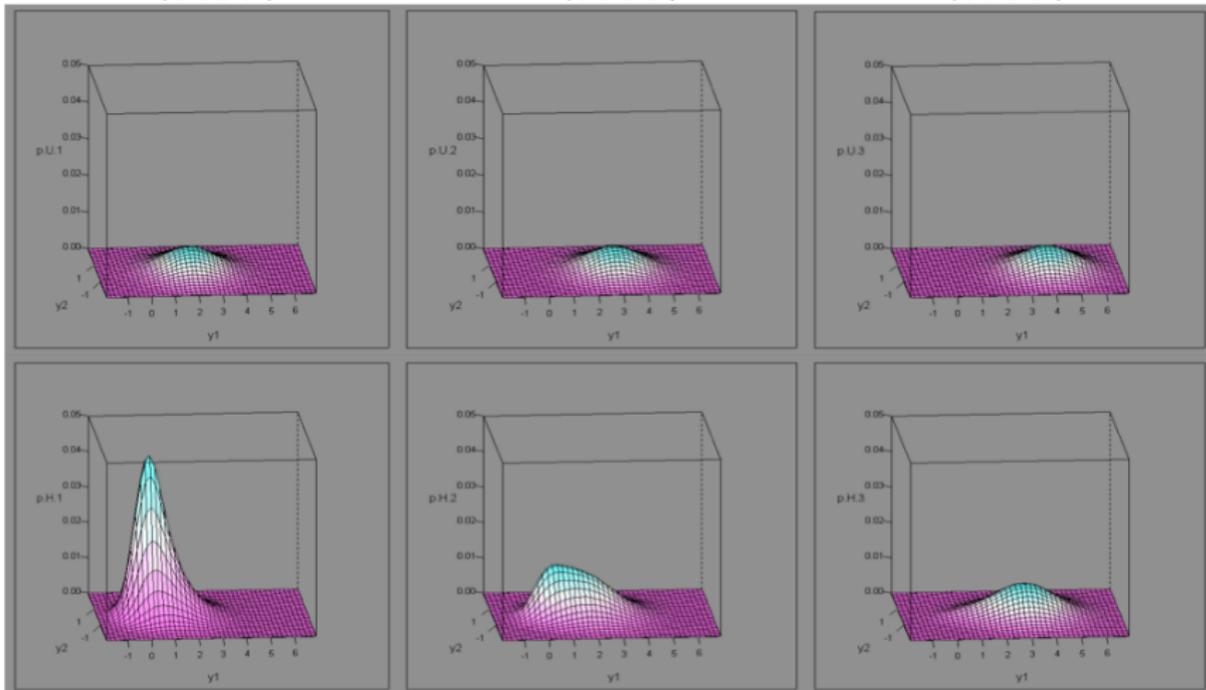

## Predictive Shrinkage in Action

- ▶ The contrast between  $p_U(y | x)$  and  $p_H(y | x)$  for various values of  $x$

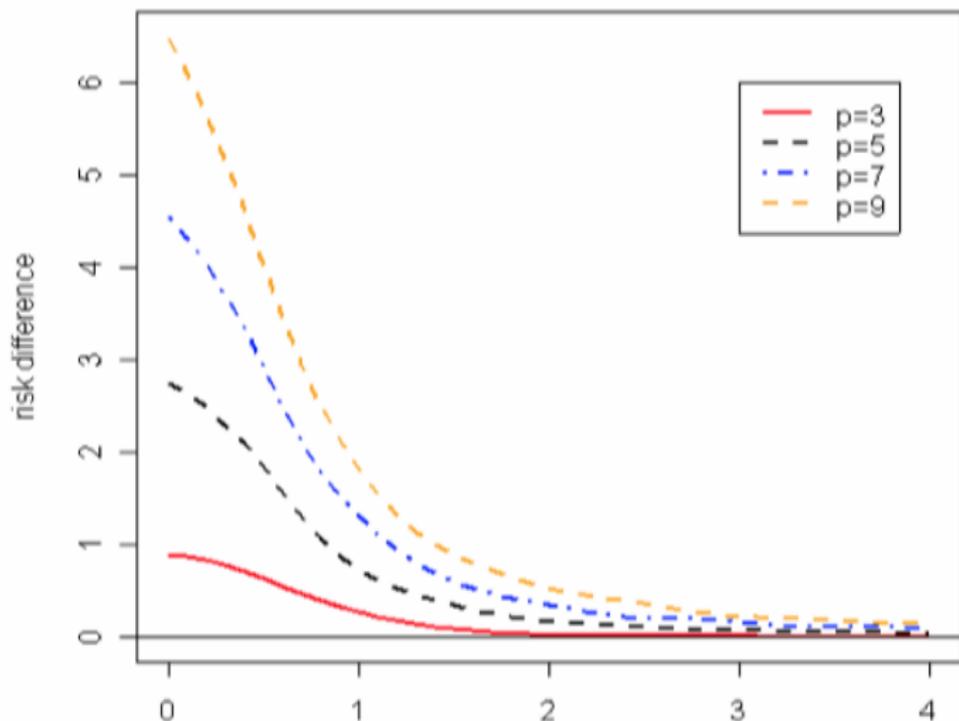
$x = (2, 0, 0, 0, 0)$

$x = (3, 0, 0, 0, 0)$

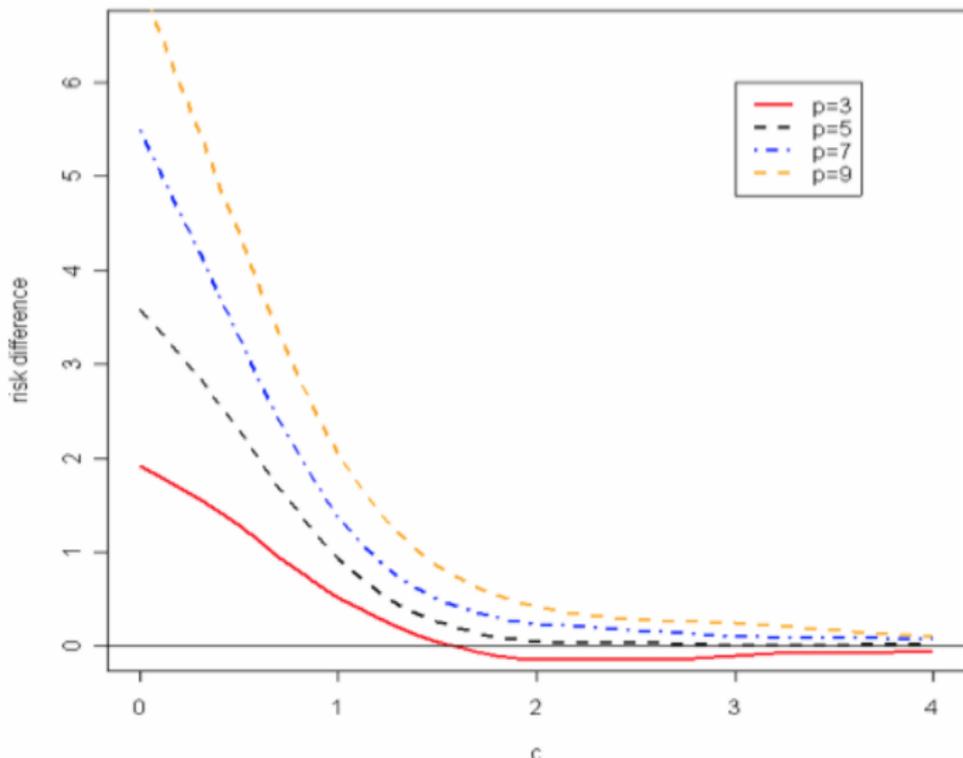
$x = (4, 0, 0, 0, 0)$



- The risk function difference  $[R_{KL}(\mu, p_U) - R_{KL}(\mu, p_H)]$  is largest at  $\mu = 0$ , and then asymptotes to 0 as  $\|\mu\| \rightarrow \infty$ .



- ▶ The risk function difference  $[R_{KL}(\mu, p_U) - R_{KL}(\mu, p_a)]$  is largest at  $\mu = 0$ , and then asymptotes to 0 as  $\|\mu\| \rightarrow \infty$ .



# Predictive Shrinkage Towards Points or Subspaces

## Predictive Shrinkage Towards Points or Subspaces

- ▶ We can trivially modify the previous priors and predictive distributions to shrink towards an arbitrary point  $b \in R^p$ .

## Predictive Shrinkage Towards Points or Subspaces

- ▶ We can trivially modify the previous priors and predictive distributions to shrink towards an arbitrary point  $b \in R^p$ .
- ▶ Consider the recentered prior

$$\pi^b(\mu) = \pi(\mu - b)$$

and corresponding recentered marginal

$$m_{\pi}^b(z; \nu) = m_{\pi}(z - b; \nu).$$

## Predictive Shrinkage Towards Points or Subspaces

- ▶ We can trivially modify the previous priors and predictive distributions to shrink towards an arbitrary point  $b \in R^p$ .
- ▶ Consider the recentered prior

$$\pi^b(\mu) = \pi(\mu - b)$$

and corresponding recentered marginal

$$m_{\pi}^b(z; \nu) = m_{\pi}(z - b; \nu).$$

- ▶ This yields a predictive distribution

$$p_{\pi}^b(y | x) = \frac{m_{\pi}^b(w; \nu_w)}{m_{\pi}^b(x; \nu_x)} p_U(y | x)$$



that now shrinks  $p_U(y | x)$  towards  $b$  rather than 0.



- ▶ More generally, we can shrink  $p_U(y | x)$  towards any subspace  $B$  of  $R^p$  whenever  $\pi$ , and hence  $m_\pi$ , is spherically symmetric.

- ▶ More generally, we can shrink  $p_U(y | x)$  towards any subspace  $B$  of  $R^p$  whenever  $\pi$ , and hence  $m_\pi$ , is spherically symmetric.
- ▶ Letting  $P_B z$  be the projection of  $z$  onto  $B$ , shrinkage towards  $B$  is obtained by using the recentered prior

$$\pi^B(\mu) = \pi(\mu - P_B \mu)$$

which yields the recentered marginal

$$m_\pi^B(z; v) := m_\pi(z - P_B z; v).$$

- ▶ More generally, we can shrink  $p_U(y | x)$  towards any subspace  $B$  of  $R^p$  whenever  $\pi$ , and hence  $m_\pi$ , is spherically symmetric.
- ▶ Letting  $P_B z$  be the projection of  $z$  onto  $B$ , shrinkage towards  $B$  is obtained by using the recentered prior

$$\pi^B(\mu) = \pi(\mu - P_B \mu)$$

which yields the recentered marginal

$$m_\pi^B(z; v) := m_\pi(z - P_B z; v).$$

- ▶ This modification yields a predictive distribution

$$p_\pi^B(y | x) = \frac{m_\pi^B(w; v_w)}{m_\pi^B(x; v_x)} p_U(y | x)$$

that now shrinks  $p_U(y | x)$  towards  $B$ .



- ▶ More generally, we can shrink  $p_U(y | x)$  towards any subspace  $B$  of  $R^p$  whenever  $\pi$ , and hence  $m_\pi$ , is spherically symmetric.
- ▶ Letting  $P_B z$  be the projection of  $z$  onto  $B$ , shrinkage towards  $B$  is obtained by using the recentered prior

$$\pi^B(\mu) = \pi(\mu - P_B \mu)$$

which yields the recentered marginal

$$m_\pi^B(z; v) := m_\pi(z - P_B z; v).$$

- ▶ This modification yields a predictive distribution

$$p_\pi^B(y | x) = \frac{m_\pi^B(w; v_w)}{m_\pi^B(x; v_x)} p_U(y | x)$$

that now shrinks  $p_U(y | x)$  towards  $B$ .



- ▶ If  $m_\pi^B(z; v)$  satisfies any of our superharmonic conditions for minimaxity, then  $p_\pi^B(y | x)$  will dominate  $p_U(y | x)$  and be minimax.

# Minimax Multiple Predictive Shrinkage

## Minimax Multiple Predictive Shrinkage

- ▶ For any spherically symmetric prior, a set of subspaces  $B_1, \dots, B_N$ , and corresponding probabilities  $w_1, \dots, w_N$ , consider the recentered mixture prior

$$\pi_*(\mu) = \sum_{i=1}^N w_i \pi^{B_i}(\mu),$$

and corresponding recentered mixture marginal

$$m_*(z; \nu) = \sum_{i=1}^N w_i m_{\pi}^{B_i}(z; \nu).$$

## Minimax Multiple Predictive Shrinkage

- ▶ For any spherically symmetric prior, a set of subspaces  $B_1, \dots, B_N$ , and corresponding probabilities  $w_1, \dots, w_N$ , consider the recentered mixture prior

$$\pi_*(\mu) = \sum_{i=1}^N w_i \pi^{B_i}(\mu),$$

and corresponding recentered mixture marginal

$$m_*(z; \nu) = \sum_{i=1}^N w_i m_{\pi}^{B_i}(z; \nu).$$

- ▶ Applying the  $\hat{\mu}_{\pi}(X) = X + \nabla \log m_{\pi}(X)$  construction with  $m_*(X; \bar{\nu})$  yields minimax multiple shrinkage estimators of  $\mu$ . (George 1986)

- ▶ Applying the predictive construction with  $m_*(z; v)$  yields

$$p_*(y | x) = \sum_{i=1}^N p(B_i | x) p_{\pi}^{B_i}(y | x)$$

where  $p_{\pi}^{B_i}(y | x)$  is a single target predictive distribution and

$$p(B_i | x) = \frac{w_i m_{\pi}^{B_i}(x; v_x)}{\sum_{i=1}^N w_i m_{\pi}^{B_i}(x; v_x)}$$

is the posterior weight on the  $i$ th prior component.

- ▶ Applying the predictive construction with  $m_*(z; v)$  yields

$$p_*(y | x) = \sum_{i=1}^N p(B_i | x) p_{\pi}^{B_i}(y | x)$$

where  $p_{\pi}^{B_i}(y | x)$  is a single target predictive distribution and

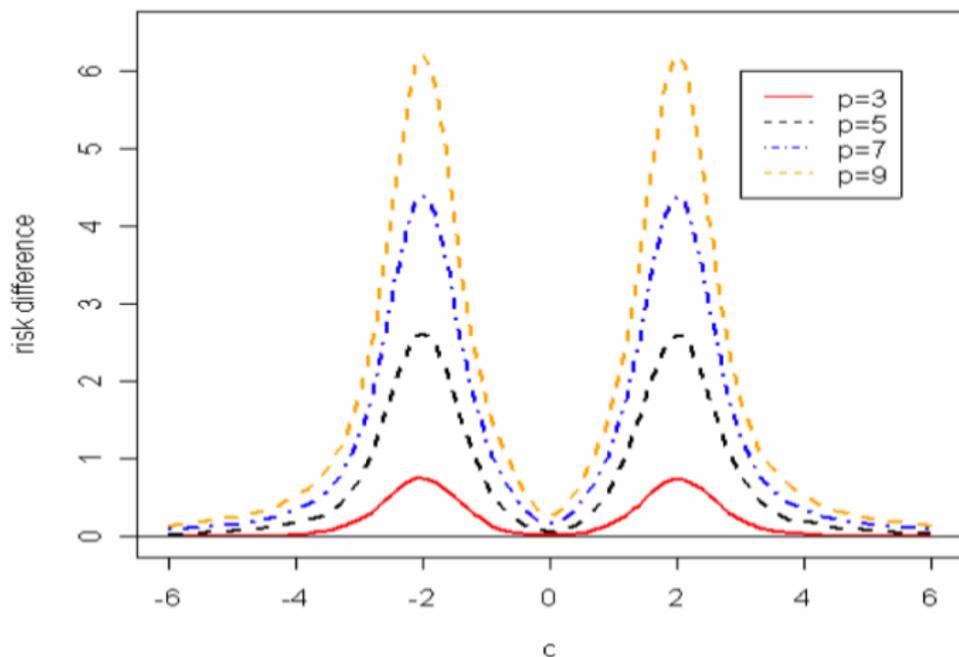
$$p(B_i | x) = \frac{w_i m_{\pi}^{B_i}(x; v_x)}{\sum_{i=1}^N w_i m_{\pi}^{B_i}(x; v_x)}$$

is the posterior weight on the  $i$ th prior component.

- ▶ **Theorem:** If each  $m_{\pi}^{B_i}(z; v)$  is superharmonic, then  $p_*(y | x)$  will dominate  $p_U(y | x)$  and will be minimax.



- ▶ The risk reduction obtained by the multiple shrinkage predictor  $p_{H^*}$  which adaptively shrinks  $p_U(y | x)$  towards the closer of the two points  $b_1 = (2, \dots, 2)$  and  $b_2 = (-2, \dots, -2)$  using equal weights  $w_1 = w_2 = 0.5$



## IV. Connecting the Estimation and Prediction Problems

## IV. Connecting the Estimation and Prediction Problems

- ▶ Comparing Stein's unbiased quadratic risk expression with our unbiased KL risk expression reveals

$$R_Q(\mu, \hat{\mu}_U) - R_Q(\mu, \hat{\mu}_\pi) = -2 \left[ \frac{\partial}{\partial \nu} E_{\mu, \nu} \log m_\pi(Z; \nu) \right]_{\nu=1}$$

## IV. Connecting the Estimation and Prediction Problems

- ▶ Comparing Stein's unbiased quadratic risk expression with our unbiased KL risk expression reveals

$$R_Q(\mu, \hat{\mu}_U) - R_Q(\mu, \hat{\mu}_\pi) = -2 \left[ \frac{\partial}{\partial \nu} E_{\mu, \nu} \log m_\pi(Z; \nu) \right]_{\nu=1}$$

- ▶ Combined with our previous KL risk difference expression reveals a fascinating connection

$$R_{KL}(\mu, p_U) - R_{KL}(\mu, p_\pi) = \frac{1}{2} \int_{\nu_w}^{\nu_x} \frac{1}{\nu^2} [R_Q(\mu, \hat{\mu}_U) - R_Q(\mu, \hat{\mu}_\pi)]_\nu d\nu$$

## IV. Connecting the Estimation and Prediction Problems

- ▶ Comparing Stein's unbiased quadratic risk expression with our unbiased KL risk expression reveals

$$R_Q(\mu, \hat{\mu}_U) - R_Q(\mu, \hat{\mu}_\pi) = -2 \left[ \frac{\partial}{\partial \nu} E_{\mu, \nu} \log m_\pi(Z; \nu) \right]_{\nu=1}$$

- ▶ Combined with our previous KL risk difference expression reveals a fascinating connection

$$R_{KL}(\mu, p_U) - R_{KL}(\mu, p_\pi) = \frac{1}{2} \int_{\nu_w}^{\nu_x} \frac{1}{\nu^2} [R_Q(\mu, \hat{\mu}_U) - R_Q(\mu, \hat{\mu}_\pi)]_\nu d\nu$$

- ▶ Ultimately it is this connection that yields the similar conditions for minimaxity and domination in both problems. Can we go further?



# Sufficient Conditions for Admissibility

## Sufficient Conditions for Admissibility

- ▶ Let  $B_{KL}(\pi, q) \equiv E_{\pi}[R_{KL}(\mu, q)]$  be the average KL risk of  $q(y | x)$  under  $\pi$ .

## Sufficient Conditions for Admissibility

- ▶ Let  $B_{KL}(\pi, q) \equiv E_{\pi}[R_{KL}(\mu, q)]$  be the average KL risk of  $q(y | x)$  under  $\pi$ .
- ▶ **Theorem** (Blyth's Method): If there is a sequence of finite nonnegative measures satisfying  $\pi_n(\{\mu : \|\mu\| \leq 1\}) \geq 1$  such that

$$B_{KL}(\pi_n, q) - B_{KL}(\pi_n, p_{\pi_n}) \rightarrow 0$$

then  $q$  is admissible.

## Sufficient Conditions for Admissibility

- ▶ Let  $B_{KL}(\pi, q) \equiv E_{\pi}[R_{KL}(\mu, q)]$  be the average KL risk of  $q(y | x)$  under  $\pi$ .
- ▶ **Theorem** (Blyth's Method): If there is a sequence of finite nonnegative measures satisfying  $\pi_n(\{\mu : \|\mu\| \leq 1\}) \geq 1$  such that

$$B_{KL}(\pi_n, q) - B_{KL}(\pi_n, p_{\pi_n}) \rightarrow 0$$

then  $q$  is admissible.

- ▶ **Theorem:** For any two Bayes rules  $p_{\pi}$  and  $p_{\pi_n}$

$$B_{KL}(\pi_n, p_{\pi}) - B_{KL}(\pi_n, p_{\pi_n}) = \frac{1}{2} \int_{v_w}^{v_x} \frac{1}{v^2} [B_Q(\pi_n, \hat{\mu}_{\pi}) - B_Q(\pi_n, \hat{\mu}_{\pi_n})]_v dv$$

where  $B_Q(\pi, \hat{\mu})$  is the average quadratic risk of  $\hat{\mu}$  under  $\pi$ .

## Sufficient Conditions for Admissibility

- ▶ Let  $B_{KL}(\pi, q) \equiv E_{\pi}[R_{KL}(\mu, q)]$  be the average KL risk of  $q(y | x)$  under  $\pi$ .
- ▶ **Theorem** (Blyth's Method): If there is a sequence of finite nonnegative measures satisfying  $\pi_n(\{\mu : \|\mu\| \leq 1\}) \geq 1$  such that

$$B_{KL}(\pi_n, q) - B_{KL}(\pi_n, p_{\pi_n}) \rightarrow 0$$

then  $q$  is admissible.

- ▶ **Theorem:** For any two Bayes rules  $p_{\pi}$  and  $p_{\pi_n}$

$$B_{KL}(\pi_n, p_{\pi}) - B_{KL}(\pi_n, p_{\pi_n}) = \frac{1}{2} \int_{v_w}^{v_x} \frac{1}{v^2} [B_Q(\pi_n, \hat{\mu}_{\pi}) - B_Q(\pi_n, \hat{\mu}_{\pi_n})]_v dv$$

where  $B_Q(\pi, \hat{\mu})$  is the average quadratic risk of  $\hat{\mu}$  under  $\pi$ .

- ▶ Using the explicit construction of  $\pi_n(\mu)$  from Brown and Hwang (1984), we obtain tail behavior conditions that prove admissibility of  $p_U(y | x)$  when  $p \leq 2$ , and admissibility of  $p_H(y | x)$  when  $p \geq 3$ .

# A Complete Class Theorem

## A Complete Class Theorem

- ▶ **Theorem:** In the KL risk problem, all the admissible procedures are Bayes or formal Bayes procedures.

## A Complete Class Theorem

- ▶ **Theorem:** In the KL risk problem, all the admissible procedures are Bayes or formal Bayes procedures.
- ▶ Our proof uses the weak\* topology from  $L^\infty$  to  $L^1$  to define convergence on the action space which is the set of all proper densities on  $R^P$ .

## A Complete Class Theorem

- ▶ **Theorem:** In the KL risk problem, all the admissible procedures are Bayes or formal Bayes procedures.
- ▶ Our proof uses the weak\* topology from  $L^\infty$  to  $L^1$  to define convergence on the action space which is the set of all proper densities on  $R^P$ .
- ▶ A Sketch of the Proof:
  1. All the admissible procedures are non-randomized.
  2. For any admissible procedure  $p(\cdot | x)$ , there exists a sequence of priors  $\pi_i(\mu)$  such that  $p_{\pi_i}(\cdot | x) \rightarrow p(\cdot | x)$  weak\* for a.e.  $x$ .
  3. We can find a subsequence  $\{\pi_{i''}\}$  and a limit prior  $\pi$  such that  $p_{\pi_{i''}}(\cdot | x) \rightarrow p_\pi(\cdot | x)$  weak\* for almost every  $x$ . Therefore,  $p(\cdot | x) = p_\pi(\cdot | x)$  for a.e.  $x$ , i.e.  $p(\cdot | x)$  is a Bayes or a formal Bayes rule.

## V. Predictive Estimation for Linear Regression

## V. Predictive Estimation for Linear Regression

- ▶ Observe  $X_{m \times 1} = A_{m \times p} \beta_{p \times 1} + \varepsilon_{m \times 1}$   
and predict  $Y_{n \times 1} = B_{n \times p} \beta_{p \times 1} + \tau_{n \times 1}$ 
  - ▶  $\varepsilon \sim N_m(0, I_m)$  is independent of  $\tau \sim N_n(0, I_n)$
  - ▶  $\text{rank}(A'A) = p$

## V. Predictive Estimation for Linear Regression

- ▶ Observe  $X_{m \times 1} = A_{m \times p} \beta_{p \times 1} + \varepsilon_{m \times 1}$   
and predict  $Y_{n \times 1} = B_{n \times p} \beta_{p \times 1} + \tau_{n \times 1}$ 
  - ▶  $\varepsilon \sim N_m(0, I_m)$  is independent of  $\tau \sim N_n(0, I_n)$
  - ▶  $\text{rank}(A'A) = p$

- ▶ Given a prior  $\pi$  on  $\beta$ , the Bayes procedure  $p_\pi^L(y | x)$  is

$$p_\pi^L(y | x) = \frac{\int p(x | A\beta)p(y | B\beta)\pi(\beta)d\beta}{\int p(x | A\beta)\pi(\beta)d\beta}$$

## V. Predictive Estimation for Linear Regression

- ▶ Observe  $X_{m \times 1} = A_{m \times p} \beta_{p \times 1} + \varepsilon_{m \times 1}$   
and predict  $Y_{n \times 1} = B_{n \times p} \beta_{p \times 1} + \tau_{n \times 1}$ 
  - ▶  $\varepsilon \sim N_m(0, I_m)$  is independent of  $\tau \sim N_n(0, I_n)$
  - ▶  $\text{rank}(A'A) = p$

- ▶ Given a prior  $\pi$  on  $\beta$ , the Bayes procedure  $p_\pi^L(y | x)$  is

$$p_\pi^L(y | x) = \frac{\int p(x | A\beta)p(y | B\beta)\pi(\beta)d\beta}{\int p(x | A\beta)\pi(\beta)d\beta}$$

- ▶ The Bayes procedure  $p_U^L(y | x)$  under the uniform prior  $\pi_U \equiv 1$  is minimax with constant risk

# The Key Marginal Representation

## The Key Marginal Representation

- ▶ For any prior  $\pi$ ,

$$p_{\pi}^L(y | x) = \frac{m_{\pi}(\hat{\beta}_{x,y}, (C'C)^{-1})}{m_{\pi}(\hat{\beta}_x, (A'A)^{-1})} p_U^L(y | x)$$

where  $C_{(m+n) \times p} = (A', B)'$  and

$$\hat{\beta}_x = (A'A)^{-1}A'x \sim N_p(\beta, (A'A)^{-1})$$

$$\hat{\beta}_{x,y} = (C'C)^{-1}C'(x', y) \sim N_p(\beta, (C'C)^{-1})$$

Risk Improvement over  $p_U^L(y | x)$

## Risk Improvement over $p_U^L(y | x)$

- ▶ Here the difference between the KL risks of  $p_U^L(y | x)$  and  $p_\pi^L(y | x)$  can be expressed as

$$R_{KL}(\beta, p_U^L) - R_{KL}(\beta, p_\pi^L) = \\ E_{\beta, (C'C)^{-1}} \log m_\pi(\hat{\beta}_{x,y}; (C'C)^{-1}) - E_{\beta, (A'A)^{-1}} \log m_\pi(\hat{\beta}_x; (A'A)^{-1})$$

## Risk Improvement over $p_U^L(y | x)$

- ▶ Here the difference between the KL risks of  $p_U^L(y | x)$  and  $p_\pi^L(y | x)$  can be expressed as

$$R_{KL}(\beta, p_U^L) - R_{KL}(\beta, p_\pi^L) = \\ E_{\beta, (C'C)^{-1}} \log m_\pi(\hat{\beta}_{x,y}; (C'C)^{-1}) - E_{\beta, (A'A)^{-1}} \log m_\pi(\hat{\beta}_x; (A'A)^{-1})$$

- ▶ Minimality of  $p_\pi^L(y | x)$  is here obtained when

$$\frac{\partial}{\partial \omega} E_{\mu, V_\omega} \log m_\pi(Z; V_\omega) < 0$$

where

$$V_\omega \equiv \omega(A'A)^{-1} + (1 - \omega)(C'C)^{-1}$$

## Risk Improvement over $p_U^L(y | x)$

- ▶ Here the difference between the KL risks of  $p_U^L(y | x)$  and  $p_\pi^L(y | x)$  can be expressed as

$$R_{KL}(\beta, p_U^L) - R_{KL}(\beta, p_\pi^L) = \\ E_{\beta, (C'C)^{-1}} \log m_\pi(\hat{\beta}_{x,y}; (C'C)^{-1}) - E_{\beta, (A'A)^{-1}} \log m_\pi(\hat{\beta}_x; (A'A)^{-1})$$

- ▶ Minimality of  $p_\pi^L(y | x)$  is here obtained when

$$\frac{\partial}{\partial \omega} E_{\mu, V_\omega} \log m_\pi(Z; V_\omega) < 0$$

where

$$V_\omega \equiv \omega(A'A)^{-1} + (1 - \omega)(C'C)^{-1}$$



- ▶ This leads to weighted superharmonic conditions on  $m_\pi$  and  $\pi$  for minimality.

## Some References

-  George, E.I., Liang, F. and Xu, X. (2006). Improved Minimax Predictive Densities Under Kullback-Leibler Loss. *Annals of Statistics*, 34 1 78–91.
-  Brown, L.D., George, E.I., and Xu, X. (2008). Admissible Predictive Density Estimation. *Annals of Statistics*, 36, 3, 1156–1170.
-  George, E.I. and Xu, X. (2008). Predictive Density Estimation for Multiple Regression. *Econometric Theory*, 24, 528–544.
-  George, E.I., Liang, F. and Xu, X. (2012). From Minimax Shrinkage Estimation to Minimax Shrinkage Prediction. *Statistical Science*, Vol. 27, No. 1 82–94.
-  Mukherjee, G. and Johnstone, I.M. (2018). Exact Minimax Estimation of the Predictive Density in Sparse Gaussian Models. *Annals of Statistics*.

Thank you!