

The James-Stein Estimator of the mean  $\mu$  of  $N_p(\mu, \Sigma)$ .

Observe  $Y_1, \dots, Y_n$  i.i.d.  $\sim N_p(\mu, \Sigma)$ . Let  $\hat{\mu}_n^1 = \bar{Y} \sim N_p(\mu, \Sigma)$ . Then:

(i)  $\hat{\mu}_n^1$  is the MLE of  $\mu$ .

(ii)  $\hat{\mu}_n^1$  is the minimum variance unbiased estimator (M.V.U.E.) of  $\mu$ :

$\hat{\mu}_n^1$  minimizes the mean square error (MSE)  $E_{\mu} \|\hat{\mu}_n^* - \mu\|^2$

among all unbiased estimators  $\hat{\mu}_n^*$  of  $\mu$ . ( $\bar{Y}$  is complete & sufficient. Apply Lehmann-Scheffe)

(iii)  $\hat{\mu}_n^1$  is the best (MSE) translation-equivariant estimator of  $\mu$ :

$$\hat{\mu}_n^*(\bar{Y} + a) = \hat{\mu}_n^*(\bar{Y}) + a \quad \forall a \in \mathbb{R}^p. \quad \text{[equivariant]}$$

(iv) For  $p=1, 2$ ,  $\hat{\mu}_n^1$  is an admissible estimator of  $\mu$ : no other estimator (linear or not) has everywhere better MSE than  $\hat{\mu}_n^1$ .

However: Stein (1956) showed that for  $p \geq 3$ ,  $\hat{\mu}_n^1$  is inadmissible.

James & Stein (1962) produced a (biased, non-linear) est.  $\hat{\mu}_n^J$  whose MSE is everywhere less than the MSE of  $\hat{\mu}_n^1$ , and substantially less for  $\mu$  near  $0$ . [Question: but is  $0$  special?] (the difference in MSE's  $\rightarrow 0$  as  $\|\mu\| \rightarrow \infty$ .)

We shall consider only the special case where  $\Sigma$  is known (for simplicity) but similar results hold if  $\Sigma$  is unknown. Thus we can assume that  $n=1$ :

Observe (\*)  $Y \sim N_p(\mu, \sigma^2 I_p)$ ,  $\sigma^2$  known. Then  $\hat{\mu}_n^1 = Y$ .

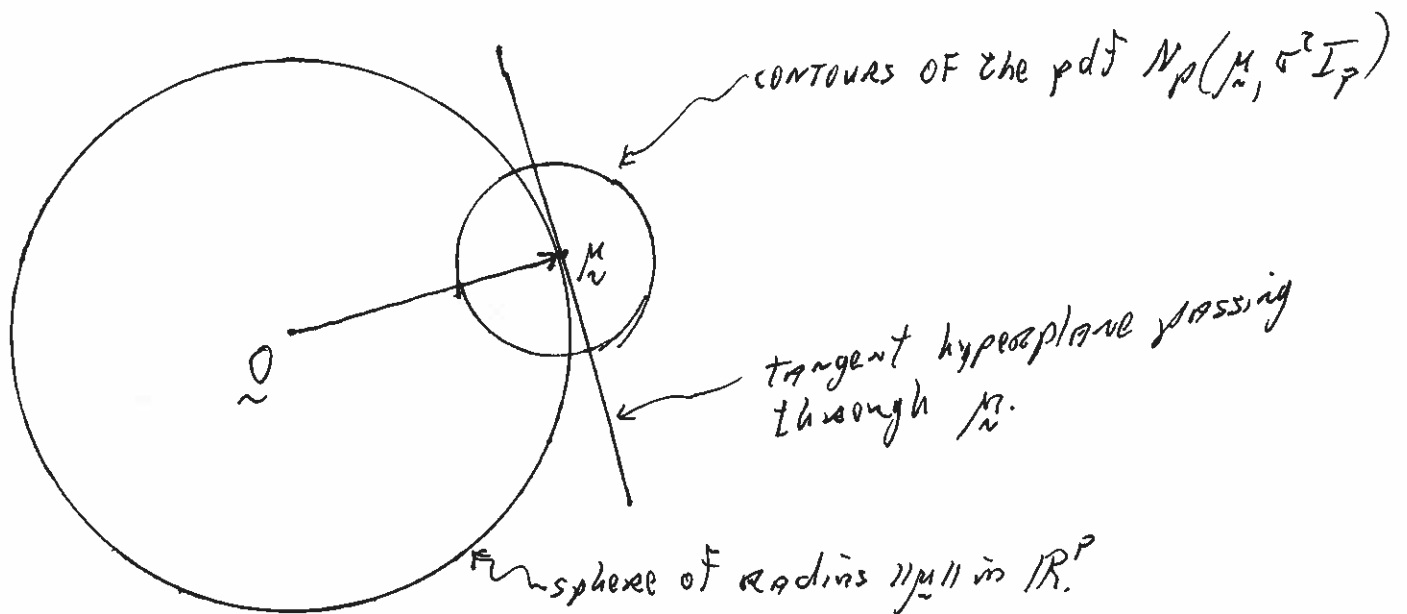
The MSE of  $\hat{\mu}_n^1$  is

$$(i) \quad E_{\mu} \|\hat{\mu}_n^1 - \mu\|^2 = E_{\mu} \|Y - \mu\|^2 = E(\sigma^2 \chi_p^2) = \sigma^2 p.$$

The JS estimator has the form  $\hat{\mu}_{JS}^1 = a(\|Y\|) \underline{Y}$ , where  $a(\|Y\|)$  is a scalar-valued function of  $\|Y\|$ . Note that  $\hat{\mu}_{JS}^1$  is not translation-equivariant, but it is orthogonally equivariant (as is  $\hat{\mu}^1$ ):

$$(2) \quad \hat{\mu}_{JS}^1(\Gamma Y) = \Gamma \hat{\mu}_{JS}^1(Y) \quad \forall \Gamma: p \times p \text{ orthogonal. ["equivariant"]}$$

The multiplier  $a(\|Y\|)$  has the form  $1 - \frac{d}{\|Y\|^2}$ , hence "shrinks"  $\underline{Y}$  toward  $\underline{0}$  (Sometimes, it shrinks  $\underline{Y}$  beyond  $\underline{0}$ , so  $\max\{a(\|Y\|), 0\}$  is preferable, but its MSE is hard to calculate analytically.) Here is a simplified version of Stein's heuristic explanation of why shrinking  $\hat{\mu}^1 \equiv \underline{Y}$  toward  $\underline{0}$  might reduce the MSE:



By symmetry, the tangent hyperplane divides  $\mathbb{R}^p$  into 2 halfspaces each having probability  $\frac{1}{2}$  under  $N_p(\mu, \sigma^2 I_p)$ .

$$\therefore P_{\mu} \left[ \underline{Y} \in \text{sphere of radius } \|\underline{\mu}\| \right] < \frac{1}{2},$$

so  $P_{\mu} \left[ \|\underline{Y}\| > \|\underline{\mu}\| \right] > \frac{1}{2}$ . This implies that  $\|\underline{Y}\|$  is "too large",

so to estimate  $\underline{\mu}$  we should "shrink"  $\underline{Y}$  toward  $\underline{0}$ . [but  $\underline{0}$  is arbitrary!]  
More precisely:

$\|Y\|^2 \sim \sigma^2 \chi_p^2 (\|M\|^2/\sigma^2)$ , so  $E\|Y\|^2 = \sigma^2 [p + \frac{\|M\|^2}{\sigma^2}] = p\sigma^2 + \|M\|^2 > \|M\|^2$ .  
 Thus  $\|Y\|^2$  is too large an estimate of  $\|M\|^2$  - in fact,  $\|Y\|^2 - p\sigma^2$  is an unbiased estimator of  $\|M\|^2$ .

The JSE can be motivated by an "empirical Bayesian" argument. Suppose we make the Bayesian assumption that  $M$  is also random, with prior distribution

(3)  $M \sim N_p(0, \lambda^2 I_p)$ . [We could replace 0 by any  $\mu_0$ ]

Then we should rewrite (\*) as  $Y|M \sim N(M, \sigma^2 I_p)$ . From these two facts, (the marginal dist'n of  $M$  and the conditional dist'n of  $Y|M$ ), we can deduce that the joint distribution of  $(Y, M)$  is the following:

(4)  $\begin{pmatrix} Y \\ M \end{pmatrix} \sim N_{2p} \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} (\sigma^2 + \lambda^2) I_p & \lambda^2 I_p \\ \lambda^2 I_p & \lambda^2 I_p \end{pmatrix} \right]$ . [Exercise]

Therefore the posterior distribution of  $M|Y$  (i.e., the conditional dist'n) is

(5)  $M|Y \sim N_p \left[ \left( \frac{\lambda^2}{\sigma^2 + \lambda^2} \right) Y, \left( \frac{\sigma^2 \lambda^2}{\sigma^2 + \lambda^2} \right) I_p \right]$  [Exercise]

The usual Bayes estimator (w.r.t. quadratic loss) is the posterior mean.

(6)  $\hat{M}_\lambda \equiv E[M|Y] = \left( \frac{\lambda^2}{\sigma^2 + \lambda^2} \right) Y$ .

This is a linear estimator of the form  $aY$  with  $a = \frac{\lambda^2}{\sigma^2 + \lambda^2} < 1$ , so  $E\|\hat{M}_\lambda - M\|^2 = E\|aY - M\|^2 = a^2 E\|Y - M\|^2 + (a-1)^2 \|M\|^2 \rightarrow \infty$  as  $\|M\| \rightarrow \infty$ , hence  $\hat{M}_\lambda$  cannot dominate the MLE  $\hat{M} \equiv Y$ . [The MSE of  $\hat{M}_\lambda$  is less than that of  $\hat{M}$  for  $M$  near 0, i.e., for  $M$  near the mean of the assumed prior dist'n (3).]

If the prior variance  $\lambda^2$  is unknown, however, then we cannot use  $\hat{M}_\lambda$ . This raises the question: can we estimate  $\lambda$ , based on  $Y$ ? If so, then we can replace  $\hat{M}_\lambda$  by  $\hat{M}_{\hat{\lambda}}$ , where  $\hat{\lambda}$  is the estimated value. The estimator  $\hat{M}_{\hat{\lambda}}$  is called an "empirical Bayes estimator".

because the prior distribution is estimated from the data.

To obtain  $\hat{\lambda}^2$ , note from (4) that the marginal dist'n of  $Y$  is

$$(7) \quad \underline{Y} \sim N_p \left[ \underline{0}, (\sigma^2 + \lambda^2) I_p \right],$$

so

$$(8) \quad \|\underline{Y}\|^2 \sim (\sigma^2 + \lambda^2) \chi_p^2.$$

Now, we can rewrite  $\hat{M}_\lambda$  as

$$(9) \quad \hat{M}_\lambda = \left( 1 - \frac{\sigma^2}{\sigma^2 + \lambda^2} \right) \underline{Y},$$

so it will suffice to estimate  $\frac{1}{\sigma^2 + \lambda^2}$ . But

$$(10), \quad \frac{1}{\|\underline{Y}\|^2} \sim \left( \frac{1}{\sigma^2 + \lambda^2} \right) \frac{1}{\chi_p^2},$$

so

$$(11) \quad E \left( \frac{1}{\|\underline{Y}\|^2} \right) = \left( \frac{1}{\sigma^2 + \lambda^2} \right) E \left( \frac{1}{\chi_p^2} \right) = \left( \frac{1}{\sigma^2 + \lambda^2} \right) \left( \frac{1}{p-2} \right)$$

provided  $p \geq 3$ .

Thus, when  $p \geq 3$ , we may estimate  $\left( \frac{1}{\sigma^2 + \lambda^2} \right)$  by  $\frac{p-2}{\|\underline{Y}\|^2}$ ; if we replace  $\frac{1}{\sigma^2 + \lambda^2}$  by this estimate in (9), we obtain the James-Stein estimator

$$(12) \quad \hat{M}_{JS} = \left[ 1 - \frac{(p-2)\sigma^2}{\|\underline{Y}\|^2} \right] \underline{Y}.$$

{When  $p=1$  or  $2$ ,  $E\left(\frac{1}{\chi_p^2}\right) = \infty$  so this motivation of (12) fails. The estimator on the right of (12) is still well defined, but it will not dominate}

THEOREM:  $E_{\underline{M}} \|\hat{M}_{JS} - \underline{M}\|^2 < E_{\underline{M}} \|\hat{M} - \underline{M}\|^2 = p\sigma^2$  For every  $\underline{M} \in \mathbb{R}^p$  ( $p \geq 3$ )

PROOF: We will prove a more general result. Define

$$(13) \quad \hat{M}_c = \left[ 1 - \frac{c\sigma^2}{\|\underline{Y}\|^2} \right] \underline{Y}, \quad [c \text{ a constant}].$$

We shall show that for every  $0 < c < 2(p-2)$ ,

$$(14) \quad E_{\underline{M}} \|\hat{M}_c - \underline{M}\|^2 < p\sigma^2 \text{ for every } \underline{M} \in \mathbb{R}^p,$$

and that the left side of (14) is minimized when  $c = p-2$ , i.e., when  $\hat{M}_c = \hat{M}_{JS}$ .

Since  $\|z\|^2 = z'z$ , we have

$$(15) \quad E_{\mu} \left\| \frac{1}{c} \tilde{M} - \tilde{M} \right\|^2 = E_{\mu} \left\| (Y - \mu) - \frac{c\sigma^2}{\|Y\|^2} Y \right\|^2$$

$$= \underbrace{E_{\mu} \left\| Y - \mu \right\|^2}_{= p\sigma^2} - 2c\sigma^2 E_{\mu} \left[ \frac{(Y - \mu)' Y}{\|Y\|^2} \right] + c^2 \sigma^4 E_{\mu} \left[ \frac{1}{\|Y\|^2} \right]$$

Now  $Y \sim N_p(\mu, \sigma^2 I_p)$ , so

$$(16) \quad E_{\mu} \left[ \frac{(Y - \mu)' Y}{\|Y\|^2} \right] = \sum_{i=1}^p E \left[ \frac{(Y_i - \mu_i) Y_i}{\|Y\|^2} \right] = \sum_{i=1}^p \int_{\mathbb{R}^p} \frac{(y_i - \mu_i) y_i}{\|y\|^2} \left[ \prod_{j=1}^p f(y_j - \mu_j) dy_j \right]$$

where

$$(17) \quad f(y_j - \mu_j) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_j - \mu_j)^2}{2\sigma^2}} \quad [\text{i.e., } Y_j \sim N_1(\mu_j, \sigma^2)]$$

Now use integration by parts:

$$(18) \quad \int_{-\infty}^{\infty} \frac{y_i}{\|y\|^2} \cdot (y_i - \mu_i) e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}} dy_i$$

$\underbrace{\hspace{10em}}_{= u} \quad \underbrace{\hspace{10em}}_{= dv}$

$$u = \frac{y_i}{y_1^2 + \dots + y_p^2}, \quad du = \frac{1}{y_1^2 + \dots + y_p^2} - \frac{2y_i^2}{(y_1^2 + \dots + y_p^2)^2}$$

$$v = -\sigma^2 e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}}$$

$$= \underbrace{\frac{y_i}{\|y\|^2} (-\sigma^2 e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}})}_{= 0} \Big|_{-\infty}^{\infty} + \sigma^2 \int_{-\infty}^{\infty} \left( \frac{1}{\|y\|^2} - \frac{2y_i^2}{\|y\|^4} \right) e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}}$$

so by (16),

$$(19) \quad E_{\mu} \left[ \frac{(Y - \mu)' Y}{\|Y\|^2} \right] = \sum_{i=1}^p \sigma^2 \int_{\mathbb{R}^p} \left( \frac{1}{\|y\|^2} - \frac{2y_i^2}{\|y\|^4} \right) \prod_{j=1}^p f(y_j - \mu_j) dy_j$$

$$= \sigma^2 E_{\mu} \left[ \frac{p}{\|Y\|^2} - \frac{2}{\|Y\|^2} \right] = \sigma^2 (p-2) E_{\mu} \left( \frac{1}{\|Y\|^2} \right)$$

Therefore, From (15),

$$(20) \quad \Delta(\mu) \equiv p\sigma^2 - E_{\mu} \left\| \frac{1}{c} \tilde{M} - \tilde{M} \right\|^2 = \sigma^4 c [2(p-2) - c] E_{\mu} \left( \frac{1}{\|Y\|^2} \right)$$

so  $\Delta(\mu) > 0$  when  $0 < c < 2(p-2)$ , and  $\Delta(\mu)$  is maximized (for all  $\mu$ ) when  $c = (p-2)$ . This completes the proof of the Theorem.

Remarks: A. When  $C = p-2$  so  $\hat{\mu}_E = \hat{\mu}_{JS}$ , the reduction in total MSE provided by  $\hat{\mu}_{JS}$  is

$$(21) \quad \Delta(\underline{\mu}) = (p-2)^2 \sigma^2 E \left[ \frac{1}{\|\underline{Y}/\sigma\|^2} \right] = (p-2)^2 \sigma^2 E \left[ \frac{1}{\chi_p^2(\delta)} \right],$$

where  $\delta = \|\underline{\mu}\|^2 / \sigma^2$ . Since the noncentral  $\chi_p^2(\delta)$ -distribution can be shown to have monotone likelihood ratio in  $\delta$ , the expectation  $E[1/\chi_p^2(\delta)]$  is decreasing in  $\delta$ , so  $\Delta(\underline{\mu})$  is decreasing in  $\|\underline{\mu}\|$ . Thus the maximum reduction in total MSE provided by  $\hat{\mu}_{JS}$  occurs when  $\underline{\mu} = 0$  and is given by

$$(22) \quad \Delta(0) = (p-2)^2 \sigma^2 E \left[ \frac{1}{\chi_p^2} \right] = (p-2) \sigma^2,$$

which is quite large for large  $p$ . On the other hand, as  $\delta \rightarrow \infty$ ,  $\chi_p^2(\delta) \rightarrow \infty$  and  $E[1/\chi_p^2(\delta)] \rightarrow 0$  [this requires verification], so

$$(23) \quad \Delta(\underline{\mu}) \rightarrow 0 \text{ as } \|\underline{\mu}\| \rightarrow \infty.$$

Thus, the total MSE of  $\hat{\mu}_{JS}$ , i.e.,  $E\|\hat{\mu}_{JS} - \underline{\mu}\|^2$ , depends on  $\|\underline{\mu}\|$  as follows:

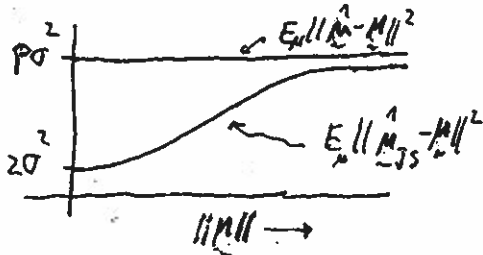


FIGURE 1.

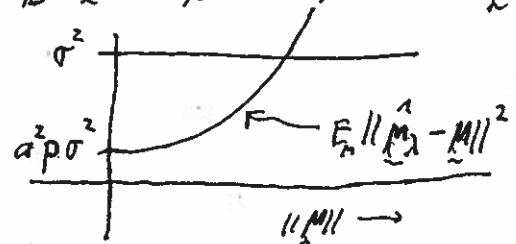


FIGURE 2.

Hence, the JS estimator  $\hat{\mu}_{JS}$  "dominates" the LSE  $\hat{\mu}$ : its total MSE is always smaller. On the other hand, the MSE of the Bayes estimator  $\hat{\mu}_\lambda$  in (6) (shown in Figure 2) exceeds that of  $\hat{\mu}$  for sufficiently large  $\|\underline{\mu}\|$ :

$$(24) \quad E\|\hat{\mu}_\lambda - \underline{\mu}\|^2 = E\|a\underline{Y} - \underline{\mu}\|^2 = E\|a(\underline{Y} - \underline{\mu}) + (a-1)\underline{\mu}\|^2 \\ = a^2 E\|\underline{Y} - \underline{\mu}\|^2 + (a-1)^2 \|\underline{\mu}\|^2 = a^2 p\sigma^2 + (a-1)^2 \|\underline{\mu}\|^2 \rightarrow \infty$$

where  $a = \frac{\lambda^2}{\sigma^2 + \lambda^2} < 1$ , so  $E\|\hat{\mu}_\lambda - \underline{\mu}\|^2 \rightarrow \infty$  as  $\|\underline{\mu}\| \rightarrow \infty$ .

B. The JS estimator  $\hat{\mu}_{JS}$  does not dominate the MLE  $\hat{\mu}$  in the strong sense

$$E_{\mu} \left( \hat{\mu} - \mu \right) \left( \hat{\mu} - \mu \right)^t - E_{\mu} \left( \hat{\mu}_{JS} - \mu \right) \left( \hat{\mu}_{JS} - \mu \right)^t = \Sigma_{MLE} - \Sigma_{JS}$$

is not positive semidefinite. The THEOREM on p. 4.4 states that  $\text{tr}(\Sigma_{MLE}) > \text{tr}(\Sigma_{JS})$ , but it can be shown that some diagonal elements of  $\Sigma_{JS}$  (i.e.,  $E_{\mu} \left( \hat{\mu}_{JS,i} - \mu_i \right)^2$ ) can exceed the corresponding diagonal elements of  $\Sigma_{MLE}$  (i.e.,  $E_{\mu} \left( \hat{\mu}_i - \mu_i \right)^2$ ).

C. The JS estimate has the form  $a(\|y\|) \cdot \underline{y}$ , where  $a(\|y\|) = \left[ 1 - \frac{(p-2)\sigma^2}{\|y\|^2} \right]$ . Since  $a(\|y\|) < 1$ ,  $\hat{\mu}_{JS}$  "shrinks"  $\underline{y}$  toward  $\underline{0}$  (mimicking the behavior of the Bayes estimator  $\hat{\mu}_2$ , which shrinks  $\underline{y}$  toward  $\underline{0} \equiv$  the prior mean). However, it can occur that  $a(\|y\|) < 0$ , (whenever  $\|y\|^2 < (p-2)\sigma^2$ ), in which case  $\hat{\mu}_{JS}$  "shrinks"  $\underline{y}$  too far. It can be shown that the modified JS estimator  $\hat{\mu}_{JS}^+ = a^+(\|y\|) \cdot \underline{y}$  dominates  $\hat{\mu}_{JS}$  (hence  $\hat{\mu}_{JS}^+$  dominates  $\hat{\mu}$ ), where  $a^+(\|y\|) = \max\{a(\|y\|), 0\}$ .

D. If  $\sigma^2$  is unknown but we have an estimator  $\hat{\sigma}^2$  of  $\sigma^2$  such that  $\hat{\sigma}^2$  is independent of  $\underline{y}$  and  $\hat{\sigma}^2 \sim \sigma^2 \chi_k^2/k$  for some  $k$ , then estimators of  $\underline{\mu}$  of the form  $\left[ 1 - \frac{c\hat{\sigma}^2}{\|y\|^2} \right] \cdot \underline{y}$  will dominate  $\hat{\mu}$  for suitable choices of  $c$ .

E. As already mentioned, if we replace the prior (3) by  $N(\underline{M}_0, \lambda^2 I_p)$ , we are led to the "JS" estimator that shrinks  $\underline{y}$  toward  $\underline{M}_0$ , namely,

$$(25) \quad \hat{\mu}_{JS}(\underline{M}_0) \equiv \left[ 1 - \frac{(p-2)\sigma^2}{\|y - \underline{M}_0\|^2} \right] \left( \underline{y} - \underline{M}_0 \right) + \underline{M}_0.$$

The choice of  $\underline{M}_0$  should not be made arbitrarily, but should be based on <sup>the</sup> prior information available. If no such info. is available, the JS estimator should not be used (this is my personal opinion).

\* F. It may be preferable to shrink toward a linear subspace, e.g. toward  $(\underline{y}, \dots, \underline{y})$ . Use p-3.