

90  
COMBINING SAMPLE INFORMATION  
IN ESTIMATING ORDERED  
NORMAL MEANS

by

Constance van Eeden and  
James V. Zidek

Technical Report #182

December 1998

# COMBINING SAMPLE INFORMATION IN ESTIMATING ORDERED NORMAL MEANS <sup>1</sup>

by

Constance van Eeden and James V Zidek  
University of British Columbia

## Abstract

In this paper we answer a question concerned with the estimation of  $\theta_1$  when  $Y_i \sim^{ind} \mathcal{N}(\theta_i, \sigma_i^2)$ ,  $i = 1, 2$ , are observed and  $\theta_1 \leq \theta_2$ . In this case  $\theta_2$  contains information about  $\theta_1$  and we show how the relevance weights in the so-called relevance weighted likelihood might be selected so that  $Y_2$  may be used together with  $Y_1$  for effective likelihood-based inference about  $\theta_1$ . Our answer to this question uses the Akaike entropy maximization criterion to find the relevance weights empirically. Although the problem of estimating  $\theta_1$  under these conditions has a long history, our estimator appears to be new. Unlike the MLE it is continuously differentiable. Unlike the Pitman estimator for this problem, but like the MLE, it has a simple form. The paper describes the derivation of our estimator, presents some of its properties and compares it with some obvious competitors. Finally, a number of open problems are presented.

*Key words:* Likelihood; maximum likelihood; weighted likelihood; estimation; admissibility; minimaxity; normal means; restricted parameter spaces; relevance weighting.

*AMS 1991 classifications:* 62F30; 62F10; 62C15; 62C20.

*STMA 1998 classifications:* 04:170; 04:010; 04:020; 04:035.

---

<sup>1</sup>Supported by a grant from the Natural Sciences and Engineering Research Council of Canada

## 1 Introduction

This article addresses broadly the problem of successfully trading off bias for precision in statistical estimation. That problem arises when an investigator has data from a population other than that of his or her inferential interest. Do these auxiliary data contain information of value for estimating parameters in the population of interest? If so, how can the bias in the auxiliary sample be traded off for precision in the required parameter estimators.

The specific problem we consider is that of estimating the mean  $\theta_1$  of a univariate normal population from which an observation  $y_1$  has been drawn. We suppose an independent observation  $y_2$  has also been drawn from another normal population with mean  $\theta_2$  when  $\theta_1 \leq \theta_2$ . Now the general questions we ask above can be stated more specifically by asking how  $y_2$  can be used in conjunction with  $y_1$  to create an estimator that improves on the estimator  $y_1$  based only on data from the first population.

Heuristics suggest an affirmative answer. The event  $y_2 < y_1$  combined with the knowledge that  $\theta_1 \leq \theta_2$  suggests  $\theta_1 \approx \theta_2$ . That suggests a better estimator of  $\theta_1$  would be obtained by taking the BLUE that would be used if the population means were equal.

We describe a new method for operationalizing these heuristics in Section 2. However, before introducing that method, we should note that a number of authors have proposed methods different from the one we obtain with our new method for exploiting  $y_2$  in the estimation of  $\theta_1$ . Unlike the classical unbiased MLE viz  $y_1$  (hereafter denoted by ULE) the alternative estimators obtained by those authors are biased like ours. However, these estimators can have substantially smaller mean-squared-errors (MSE's) than their classical counterpart over portions of the parameter space deemed to be of particular importance. At the same time, their MSE's are either smaller or not appreciably larger over the rest of the parameter space than the MSE of the ULE. Thus an effective bias-variance trade-off is indeed possible; information in the sample from the second population can help in estimating the mean of the first.

In Section 3 we describe estimators developed by other authors to make that trade-off. However before doing so, we develop in Section 2 new estimators using an extension of Fisher's classical likelihood that Hu (1994) introduces and calls the "Relevance Weighted Likelihood REWL." It generalizes the local likelihood defined in the context of non-parametric regression by Tibshirani and Hastie (1987) that was extended as a local likelihood by Staniswalis (1989) and as a quasi-local-likelihood by Fan, Heckman and Wand (1995).

In contrast to the local likelihood, the REWL can be a global likelihood and in one of the applications developed by Hu and Zidek (1997), it is shown how the celebrated James-Stein estimator can be found as a maximum (relevance weighted) likelihood estimator when the relevance weights are estimated from the data.

The relevance weights allow bias to be traded for precision in the likelihood setting, as bias is traded for variance in the non-parametric regression setting. The need for such a theory has become increasingly important as the scale of modern experimental science has grown in its space-time scales thanks to demand (eg. environmental science) combined with feasibility (eg. through information technology). On these scales, the replicated experiment seems completely unrealistic as an experimental paradigm, leading to the need for a theory that embraces bias without sacrificing the goals of efficiency and precision enshrined in Fisher's foundational works.

The theory described in Section 2 enables the bias-precision trade-off to be made without relying on the Bayesian approach (see Berger 1985). The latter permits the bias-variance trade-off to be made in a conceptually straightforward manner. Reliance on empirical Bayes methods softens the demands for realistic prior modeling in complex problems. Efron (1996) illustrates the empirical Bayes approach in such problems and uses the term "relevance" in a manner similar to that of Hu (1994).

Our theory is proposed as a simpler alternative to the empirical Bayesian approach for use in complex problems. The REWL offers such an approach and we will try to demonstrate that in this article. At the same time we gain a theory that formally links a diverse collection of statistical domains such as weighted least squares, non-parametric regression, meta-analysis and shrinkage estimation. Starting with the likelihood in these domains yields new methods and suggests new problems as we will attempt to show. At the same time, the REWL comes with an (as yet incomplete) underlying general theory including extensions of Wald's theory for the maximum likelihood estimator (Hu 1997).

In Section 3 we address study the bias-variance trade-off made by a number of biased estimators proposed as solutions to the problem central to this paper. Included is the estimator we propose in Section 2. Numerical assessments of their properties point to a number of conjectures and questions listed in that section for deeper analysis in Section 4.

In Section 4 we prove a number of the conjectures in Section 3 and at the same time answer a number of the questions posed there. However many of the conjectures remain unproven and questions unanswered.

These are listed in the concluding Section 5. There as well we summarize the results



of our inquiry and the possible value of the REWL-based methodology.

## 2 Relevance Weighted Likelihood Estimation

In this section we describe for completeness the relevance weighted likelihood in the general case and then apply it to the specific problems of interest in this paper. Assume  $\{Y_i\}$  are independently distributed random variables or vectors, each having an associated population distribution with probability density and cumulative distribution (PDF and CDF, respectively)  $f_i$  and  $F_i$ . Let  $\mathbf{Y} = (Y_1, \dots, Y_n)$  be the vector or matrix of these measurable attributes.

From each population  $i$ ,  $n_i \geq 0$  items are randomly and independently sampled, yielding  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})$ ,  $Y_{ij}$  representing the  $Y_i$  measured on the  $j$ -th item sampled from the  $i$ -th population  $j = 1, \dots, n_i$ ,  $i = 1, \dots, n$  (the null vector when  $n_i = 0$ ). Assume the  $Y_{ij}$ ,  $j = 1, \dots, n_i$  are independent as well as identically distributed, each having its associated population distribution. Denote the realization of  $\mathbf{Y}_i$  by  $\mathbf{y}_i$ ,  $i = 1, \dots, n$ .

In this paper inferential interest concerns attributes of population 1. However in general Hu and Zidek (1997) consider other possibilities such as simultaneous inference about parameters of all the populations.

Starting from the Akaike entropy maximization principle (1973, 1977, 1978, 1982, 1983, 1985), Hu and Zidek (1997) derive the REWL in the non-parametric and parametric cases. To be precise they suppose (when the  $Y$  are discrete) that a predictive distribution say  $g$  of  $Y_1$  must be chosen to maximize  $\int \log g(y) dF_1(y)$  where  $F_1$  denotes the true "conceptual" population distribution for the first population. This maximization must be done subject to knowledge that  $F_1$  resembles each of the other  $F_j$ ,  $j \neq 1$ , that is subject to  $\int \log g(y) dF_j(y) > c_j$ ,  $j \neq 1$  for specified  $\{c_j, j \neq 1\}$ . A Lagrangian argument then implies that  $g$  maximizes a linear combination the  $\int \log g(y) dF_j(y)$ ,  $j = 1, \dots, n$ . However since the  $\{F_j\}$  are unknown they are estimated by  $\{\hat{F}_j\}$  their empirical distribution functions. When only one observation  $y_j$  is available from population  $j = 1, \dots, n$ , the empirical distribution for that population becomes a point mass at that observation.

In any event, with these heuristics the optimum  $g$  maximizes the non-parametric relevance likelihood function that viewed as a function of  $g$  is

$$g \rightarrow \prod_{j=1}^n \prod_{i=1}^{n_i} g^{1/n_i}(y_{ji}). \quad (2.1)$$

Similar heuristics apply to the case of interest in this paper, i.e. the parametric case, where for the likelihood we have instead

$$\theta \rightarrow \prod_{j=1}^n \prod_{i=1}^{n_j} f_i^{\lambda_{ij}/n_j}(y_{ji} | \theta_i) \quad (2.2)$$

where  $\theta = (\theta_1, \dots, \theta_n)$ . In both cases  $\lambda_{ij} \geq 0$  and take  $\lambda_{ij}/n_j = 0$  when  $n_j = 0$  for all  $i$  and  $j$ .

The relevance weights  $\{\lambda_{ij}\}$  enable the investigator to trade off bias for precision in estimating the likelihood for population 1 using the data from the remaining populations. Ideally the choice of these weights (equivalently the specification of the  $\{c_j\}$  above) will be context dependent. However Hu and Zidek (1997) suggest a general method for their selection based on a suggestion of Stigler (1990). That method again based on the use of the maximization of entropy approach with follow-up estimation is the one used in this paper. Rather than describe it in general we demonstrate it below in specific problems.

The MREWLE for  $\theta_i$  is found by maximizing (2.2). Hu (1997) shows that the theory of Wald for the classical MLE extends to the MREWLE under a suitable adaptation of Wald's assumptions.

We apply the non-parametric REWL to the case of two normal populations  $Y_i \sim N(\theta_i, \sigma_i^2)$  for which the  $\{\sigma_i^2\}$  are known  $i = 1, 2$ . Now  $n_1 = n_2 = 1$  for the two populations involved and for simplicity we denote the relevance weights by  $\lambda_{i1} = \lambda_i$ ,  $i = 1, 2$  for those populations. The MREWLE for  $\theta_1$  or WLE for short is easily shown to be

$$\delta_{WLE}(Y_1, Y_2) = Y_1 + W\alpha$$

where  $W = Y_2 - Y_1$  and  $\alpha \in [0, 1]$  obtains from the relevance weights and needs to be specified. The relevance weight ratio defines  $\alpha$  through

$$\frac{\lambda_2}{\lambda_1} = \frac{\sigma_2^2}{\sigma_1^2} \frac{\alpha}{1 - \alpha}. \quad (2.3)$$

The maximization of entropy criterion above may be applied to find relevance weights. That criterion leads to the minimization of the MSE in this case of normal population distributions. Hence the optimal choice of  $\alpha$  if  $\Delta = \theta_2 - \theta_1$  were known would be

$$\alpha_{\text{optimal}} = \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2 + \Delta^2}.$$

However, since  $\Delta$  is unknown it must be estimated. The appropriate estimator for

the case considered in Section 3 where  $\Delta \geq 0$  would be

$$\hat{\alpha}(W) = \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2 + CW_+^2} \quad (2.4)$$

where  $W_+ = \max\{0, W\}$  and  $C$ , the “attenuation” constant must be selected by the investigator.

This approach yields a smooth estimator since  $\alpha$  is “fitted” to the data only after the MSE has been computed. In particular it is a differentiable function of  $W$  in contrast to the truncated MLE of  $\theta_1$  which is not. Now the performance of the proposed estimator needs to be explored and we do this both theoretically and numerically in the next section.

However, Hu and Zidek (1997) emphasize that the specification of the relevance weights should best be done in the context of the specific inferential context. This suggestion may be followed in the restricted means problem above since a variety of estimators that exploit  $Y_2$  in the estimation of  $\theta_1$  have already been proposed. Moreover each may be written in the form above for the MREWLE with an estimated  $\alpha$ . Thus each entails an implicit choice of the relevance weight ratio that can be exploited through the equation above relating that ratio to  $\alpha$ . In this paper we will explore these various choices and compare the associated estimators in the next section.

### 3 The Bias-Variance Trade-off.

The bias-variance trade-off goes back at least as far as Stein’s discovery that it could be made in the simultaneous estimation of independent normal population means. That celebrated discovery stimulated the study of biased estimation. The feasibility of the trade-off was demonstrated in a wide variety of contexts. One such context was that of the present paper wherein a number of biased estimators of ordered normal means were proposed.

We now examine that trade-off and the way it has been made by those estimators. Specifically we compare five estimators of  $\theta_1$  based on  $(Y_1, Y_2)$ . They are:  $\delta_{WLE}(Y_1, Y_2)$  the WLE as defined and discussed in Section 2;  $\delta_{MLE}(Y_1, Y_2)$  the MLE, i.e. the first co-ordinate of the MLE for  $(\theta_1, \theta_2)$  under the restriction  $\theta \leq \theta_2$ ;  $\delta_{ULE}(Y_1, Y_2) = Y_1$  the unrestricted MLE of  $\theta_1$  based on  $Y_1$ ;  $\delta_{MIN}(Y_1, Y_2)$  the minimum of  $Y_1$  and  $Y_2$ ; and  $\delta_P(Y_1, Y_2)$  the so-called Pitman estimator, i.e. the first co-ordinate of the generalized Bayes estimator of  $(\theta_1, \theta_2)$ , that estimator being computed from the uniform prior on



$\{(\theta_1, \theta_2) \mid \theta_1 \leq \theta_2\}$ . These estimators are given explicitly below:

$$\delta_{WLE}(Y_1, Y_2) = Y_1 + W\hat{\alpha}(W) \quad \text{where} \quad \hat{\alpha}(W) = \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2 + CW_+^2}; \quad (3.1)$$

$$\left. \begin{aligned} \delta_{MLE}(Y_1, Y_2) &= \min \left( Y_1, \frac{\sigma_1^2 Y_1 + \sigma_2^2 Y_2}{\sigma_1^2 + \sigma_2^2} \right) = Y_1 + \frac{1}{1+\tau} W_- \\ \text{with } \tau &= \sigma_2^2/\sigma_1^2 \quad \text{and} \quad W_- = \min(0, W); \end{aligned} \right\} \quad (3.2)$$

$$\delta_{VLE}(Y_1, Y_2) = Y_1; \quad (3.3)$$

$$\delta_{MIN}(Y_1, Y_2) = \min(Y_1, Y_2) = Y_1 + W_-; \quad (3.4)$$

$$\delta_P(Y_1, Y_2) = Y_1 - \sqrt{\frac{\sigma_1^4}{\sigma_1^2 + \sigma_2^2}} \frac{\phi\left(\frac{W}{\sqrt{\sigma_1^2 + \sigma_2^2}}\right)}{\Phi\left(\frac{W}{\sqrt{\sigma_1^2 + \sigma_2^2}}\right)}. \quad (3.5)$$

The “attenuation constant”  $C$  in the expression above for WLE can be adjusted to reduce dependence on  $Y_2$ . Unless otherwise noted, we will take  $C = 1$  in the ensuing discussion.

The Pitman estimator was proposed and studied by Cohen and Sackrowitz (1970). Note, however, that our formula for  $\delta_P(Y_1, Y_2)$  is not the same as the one given by Cohen and Sackrowitz. They claim, erroneously, that one can suppose, without loss of generality, that one of the two variances equals 1, making their formula valid for that special case only.

**Remark 3.1** *Note the differences in the way the above estimators depend on  $\sigma_1^2$  and  $\sigma_2^2$ . The estimators  $Y_1$  and  $\min(Y_1, Y_2)$  are independent of these variances, the relevance weighted and the Pitman estimator depend on both of them, while the MLE depends on  $\sigma_1^2$  and  $\sigma_2^2$  only through their ratio.*

We begin by examining in Figure 1 the MSE’s of the estimators plotted as functions of  $\Delta = \theta_2 - \theta_1$ .

For definiteness we have chosen  $\sigma_1 = \sigma_2 = 1$  (and  $C = 1$  in the WLE). We consider cases below where the population variances are unequal. For that reason we will in



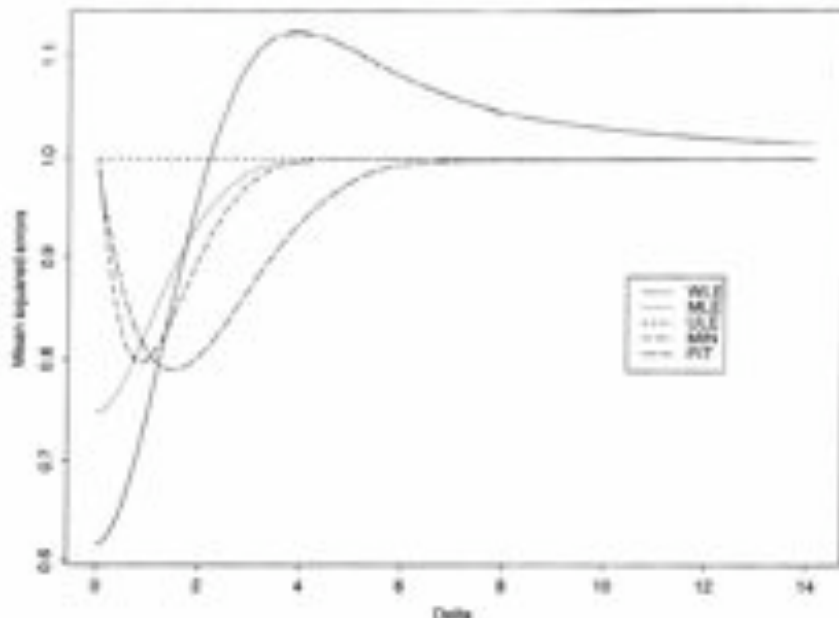


Figure 1: Graphs of the Mean Squared Errors: Selected Estimators.

general divide all the MSE's by  $\sigma_1^2$  to enable us to compare MSE plots. Therefore in all such plots the one for ULE has constant value 1 for all  $\Delta$  whatever be  $\sigma_1$ . As the classical (uniform minimum variance unbiased) estimator of  $\theta_1$ , the ULE provides a natural benchmark for assessing the performance of the alternatives considered in this paper.

The MSE of another classical estimator, the MLE also appears in Figure 1. It appears to be uniformly smaller than that of the ULE but the two are in close agreement for large  $\Delta$ . That agreement encourages optimism about the quality of the ULE since generally the MLE performs well. We express our optimism in the following conjecture.

**Conjecture 1:** ULE and MLE are minimax estimators.

At the same time the MLE appears to dominate the ULE leading us to a second conjecture.

**Conjecture 2:** The ULE is inadmissible and dominated by the MLE.

Furthermore we are led to a question:

**Question 1:** Is the MLE admissible?

Figure 1 shows the MSE for the WLE (as well as the MLE) to be much smaller than

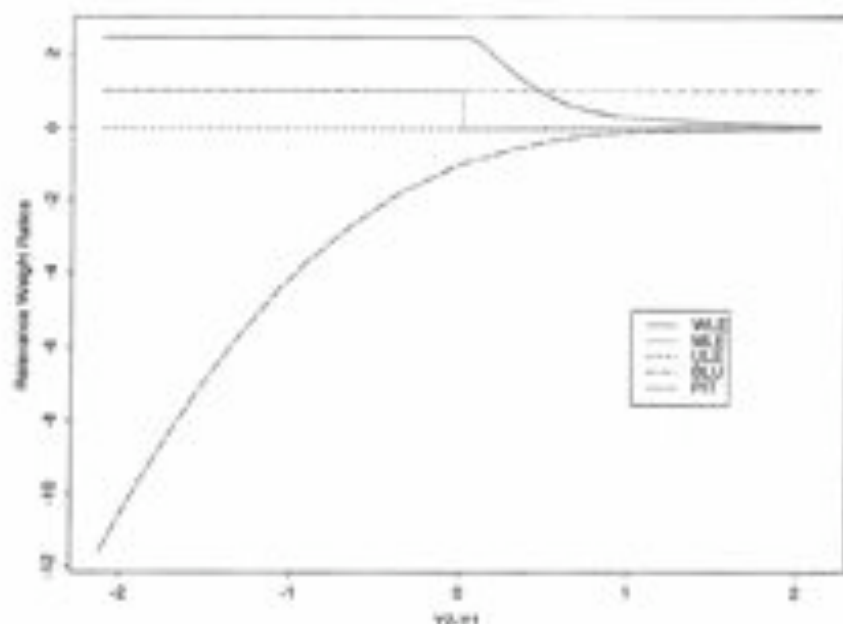


Figure 2: Graphs of Relevance Weight Ratios for Population 2 Versus 1 for Selected Estimators.

that of the ULE for small  $\Delta$ -values. Moreover its MSE resembles the MLE's for such values.

How do the MLE and the WLE achieve their seeming superiority over the ULE? The immediate answer is that they exploit the information in  $Y_2$  and they do so in a similar way. Figure 2 confirms this. That figure depicts for all estimators other than the MIN, the implied or explicit relevance weight ratios as functions of  $W = Y_2 - Y_1$ . The ratios for the MLE and WLE are broadly similar. However the WLE - ratio decreases to zero more slowly than that of the MLE. Thus it makes more liberal use of that information than does MLE. (It does so at the cost of greater bias.)

As noted above we can reduce WLE's dependence on  $Y_2$  by increasing the value of the attenuation constant. In Figure 3 we see the relevance weight ratio for the WLE approaching that of the MLE when  $C$  is chosen to be 29. Moreover, Figure 4 shows their associated MSE's to be very similar when the MLE is highly attenuated. In particular that of the WLE remains substantially smaller than that of the ULE for small  $\Delta$  values.

To gain a better understanding of how that superior performance is achieved by the WLE and the MLE relative to the other two estimators we turn to Figure 5 and see the bias functions of the various estimators. Note the comparatively small absolute

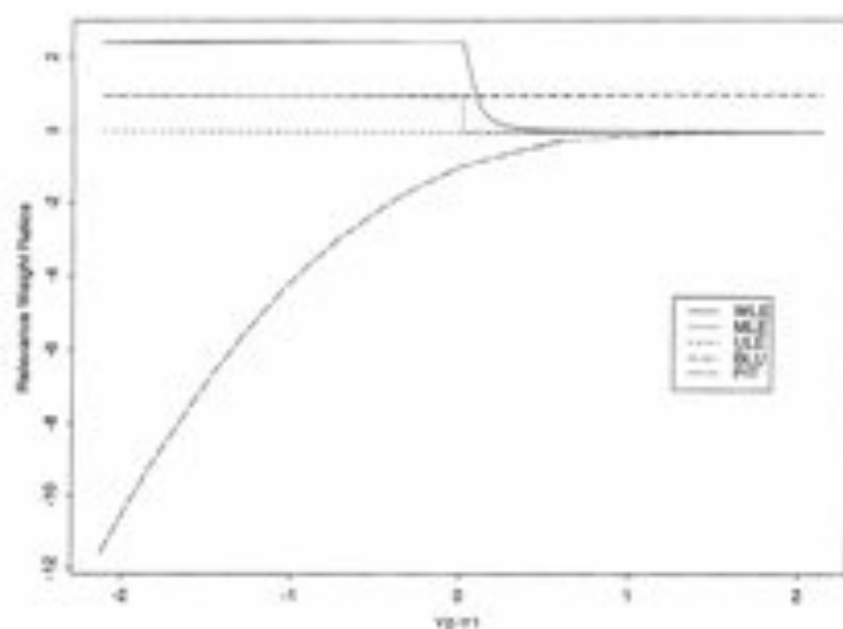


Figure 3: Graphs of Relevance Weight Ratios for Population 2 Versus 1: Selected Estimators When the WLE is Highly Attenuated, i.e  $C = 29$ .

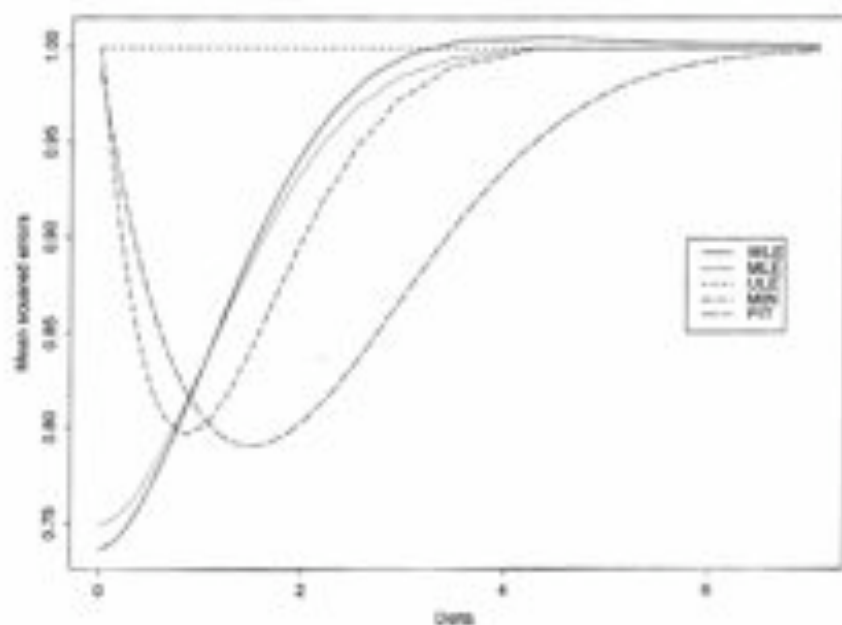


Figure 4: Graphs of the Mean Squared Errors for Selected Estimators When the WLE is Highly Attenuated, i.e.  $C = 29$ .

biases for both estimators when  $\Delta$  is close to zero compared to those of PIT and MIN. So we see that both WLE and MLE gain their superiority over ULE by aggressively exploiting the relevant information in  $Y_2$  to reduce their variances while controlling their biases for small  $\Delta$ .

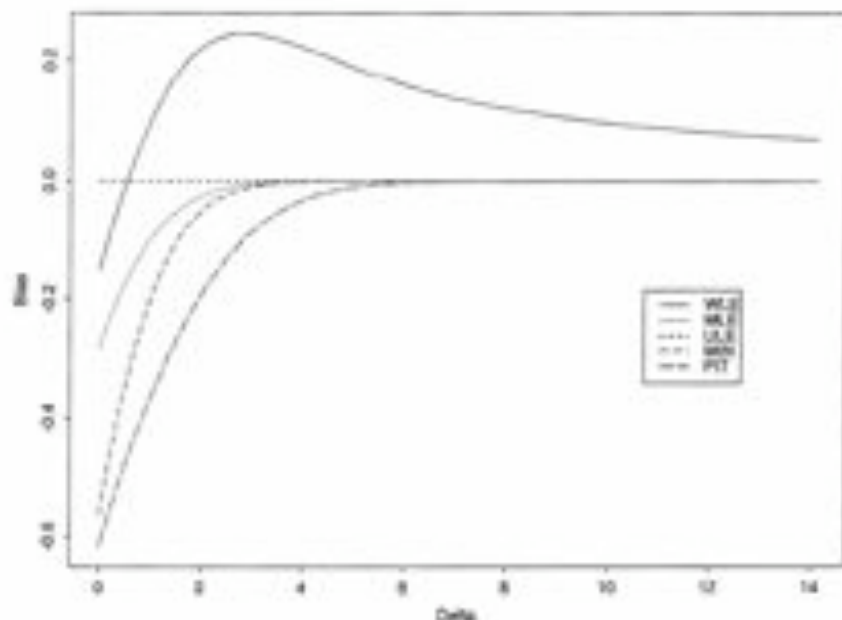


Figure 5: Graphs of Bias Functions for selected Estimators.

At the same time, Figure 1 shows that as  $\Delta$  grows larger the MSE for the WLE increases and eventually exceeds that of the ULE. Based on the earlier conjectures we make the next conjecture.

**Conjecture 3:** The WLE is not a minimax estimator.

We have the same question for the WLE as we had above for the MLE:

**Question 2:** Is the WLE admissible?

Unlike the MLE, the WLE is (twice) differentiable. The well-known necessary condition for admissibility (see Brown (1986, Theorem 4.23)) that estimators must be regular functions of the data encourages the belief that the answer to Question 2 might be “Yes”.

The remaining two estimators under consideration in this paper, PIT and MIN also seem to successfully trade bias for variance. In fact Figure 1 suggests the next conjectures.



**Conjecture 4:** PIT and MIN dominate ULE.

**Conjecture 5:** Both PIT and MIN are minimax when the population variances are identical.

That figure shows that neither estimator performs especially well when  $\Delta$  is close to zero. (They effect the bias-variance trade-off in rather subtle ways.) Nevertheless they could be admissible suggesting the next question.

**Question 3:** Are MIN and PIT admissible when the population variances are equal?

Observe in Figure 1 that the ULE-MSE uniformly exceeds that of the Pitman estimator. Moreover the comparative advantage of the Pitman estimator obtains not at  $\Delta = 0$  but rather for  $\Delta$  around 2. To interpret this observation note that the Pitman prior does not put high weight on  $\theta_1 = \theta_2$ . In fact its uniform prior on the range of  $(\theta_1, \theta_2)$  forces PIT to optimize by requiring a negative relevance weight ratio (see Figure 2). It “pushes away” the information in  $Y_2$  when the WLE and MLE embrace it (when  $\Delta = 0$ ) since under the prior this possibility would be remote. Instead PIT saves the trade-off for values of more realistic  $\Delta$ 's under the assumed prior. Nevertheless like the other alternatives to the ULE considered here other than WLE, PIT proves to be negatively biased; it tends to underestimate  $\theta_1$  (see Figure 5).

MIN succeeds in making the bias-variance tradeoff (see Figure 1) but the mechanism by which it does this proves elusive. Its weight ratio for the MIN cannot even be plotted in Figure 2, being infinite when  $W < 0$  since in that case the estimator puts all the weight on  $Y_2$  and none on  $Y_1$ . On the other hand when  $W \geq 0$  that ratio becomes zero. How does MIN so successfully exploit  $Y_2$ ? The answer seems to be that since  $\Delta \geq 0$ ,  $Y_2 \leq Y_1$  suggests  $Y_1$  is an overestimate of  $\theta_1$ . We can then profitably shrink it down to  $Y_2$ . To test this explanation we consider its implication when  $\sigma_2 < \sigma_1$  when  $Y_2$  is a measurement of higher quality than  $Y_1$  (even if biased as an estimator of  $\theta_1$ ). In this case  $Y_2$  would indicate quite reliably when  $Y_1$  overestimates  $\theta_1$ .

Figure 6 validates this heuristic reasoning. The relative gain in MIN's performance over that of ULE exceeds its gain when the population variances are unequal.

On the other hand the explanation also suggests that when  $Y_2$  is of low quality it will not help much to show when  $Y_1$  overestimates  $\theta_1$ . Again the implication is validated, this time by Figure 7. MIN now performs poorly against the other estimators as measured by its MSE.

These numerical assessments thus tend to support our explanation of how MIN works and when it would perform well. It also points to the desirability of making MIN

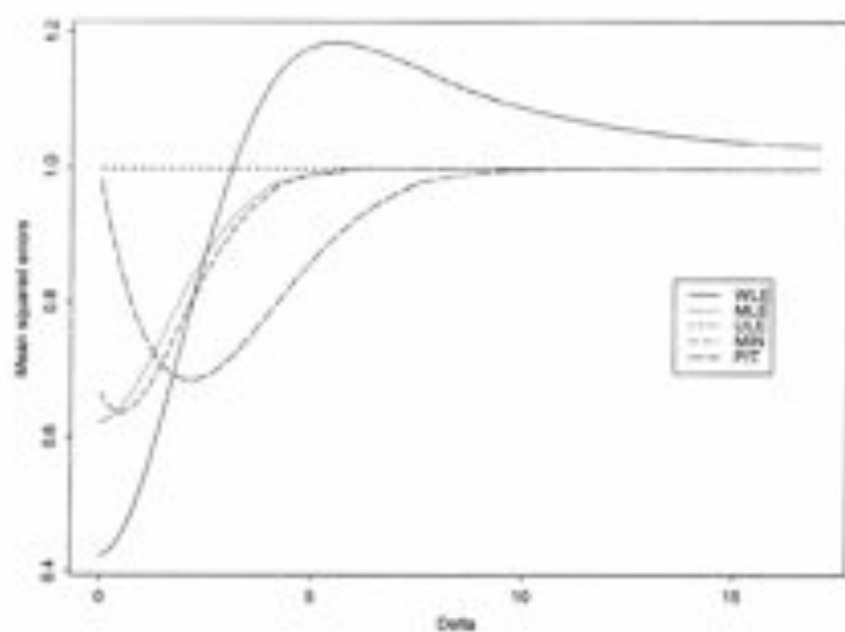


Figure 6: Graphs of Mean Squared Error Functions for Selected Estimators When Population 1 (Variance = 3) is Overdispersed Relative to 2 (Variance = 1).

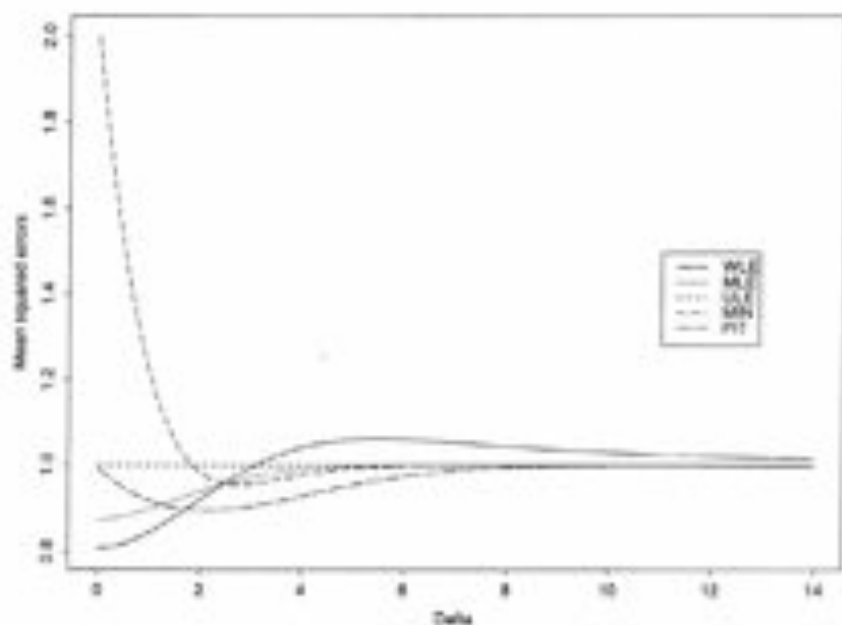


Figure 7: Graphs of Mean Squared Error Functions for Selected Estimators When Population 2 (Variance = 3) is Overdispersed Relative to 1 (Variance = 1).

depend on the population variances. This leads to our final question:

**Question 4:** Can a minimax estimator resembling MIN be found for estimating  $\theta_1$  when  $\sigma_2 > \sigma_1$ ?

## 4 Ordered Normal Means

In this section we answer some of the questions raised by the analysis of the last section. Those answers are stated as theorems whose proofs can be found in the Appendix. We begin by stating in the next theorem the mean-squared-errors of the estimators considered in the last section. There  $\Delta = \theta_2 - \theta_1$ ,  $\sigma^2 = \sigma_1^2 + \sigma_2^2$  and  $\theta = (\theta_1, \theta_2)$ .

**Theorem 4.1** For the MSEs we have :

1. The MSE of  $\delta_{WLE}$  is given by

$$\begin{aligned} \mathcal{E}_\theta(\delta_{WLE}(Y_1, Y_2) - \theta_1)^2 = \\ \sigma_1^2 + \frac{2}{1+\tau} \mathcal{E}_\theta W \dot{\alpha}(W)(\Delta - W) + \mathcal{E}_\theta \dot{\alpha}^2(W) W^2 = \\ \sigma_1^2 + \sigma^2 \sigma_1^2 \mathcal{E}_\theta \frac{W^2 (2I(W > 0) + (1+\tau)^{-1}) - 2(\sigma_1^2 + \sigma_2^2)}{(\sigma_1^2 + \sigma_2^2 + W^2)^2}; \end{aligned} \quad (4.1)$$

2. The MSE of  $\delta_{MLE}$  is given by

$$\begin{aligned} \mathcal{E}_\theta(\delta_{MLE}(Y_1, Y_2) - \theta_1)^2 = \\ \sigma_1^2 + \frac{1}{(1+\tau)^2} \left( 2\Delta \mathcal{E}_\theta W I(W < 0) - \mathcal{E}_\theta W^2 I(W < 0) \right) = \\ \sigma_1^2 + \frac{1}{(1+\tau)^2} \left\{ (\Delta^2 - \sigma^2) \left( 1 - \Phi\left(\frac{\Delta}{\sigma}\right) \right) - \Delta \sigma \phi\left(\frac{\Delta}{\sigma}\right) \right\}; \end{aligned} \quad (4.2)$$

3. The MSE of  $\delta_{VLE}$  is given by

$$\mathcal{E}_\theta(\delta_{VLE}(Y_1, Y_2) - \theta_1)^2 = \sigma_1^2; \quad (4.3)$$

4. The MSE of  $\delta_{MIN}$  is given by

$$\begin{aligned} \mathcal{E}_\theta(\delta_{MIN}(Y_1, Y_2) - \theta_1)^2 = \\ \sigma_1^2 + \frac{1}{1+\tau} \mathcal{E}_\theta W I(W < 0)(2\Delta - W(1-\tau)) = \\ \sigma_1^2 + \left( \Delta^2 - \sigma^2 \frac{1-\tau}{1+\tau} \right) \left( 1 - \Phi\left(\frac{\Delta}{\sigma}\right) \right) - \Delta \sigma \phi\left(\frac{\Delta}{\sigma}\right); \end{aligned} \quad (4.4)$$

5. The MSE of the Pitman estimator is given by

$$\mathcal{E}_\theta(\delta_P(Y_1, Y_2) - \theta_1)^2 = \sigma_1^2 - \frac{\sigma_1^4}{\sigma^2} \Delta \mathcal{E}_\theta \frac{\phi\left(\frac{W}{\sigma}\right)}{\Phi\left(\frac{W}{\sigma}\right)}. \quad (4.5)$$

In the next theorem some of the MSEs are compared. Like all comparisons between MSEs in this paper, they are made on the restricted parameter space  $\Theta = \{\theta \mid \theta_1 \leq \theta_2\}$ . Thus, an estimator  $\delta$  is inadmissible for estimating  $\theta_1$  if there exists an estimator  $\delta^*$  dominating it on  $\Theta$  and  $\delta$  is minimax if it minimizes, among estimators  $\delta^*$ ,  $\sup_{\theta \in \Theta} R(\delta^*, \theta)$ , where  $R(\delta^*, \theta)$  is the MSE of  $\delta^*$  at the parameter point  $\theta$ .

**Theorem 4.2** Each of the estimators  $\delta_{MLE}$  and  $\delta_P$  dominates  $\delta_{VLE}$ . The estimator  $\delta_{MIN}$  dominates  $\delta_{VLE}$  when  $\tau \leq 1$ . The MSE's of  $\delta_{MIN}$  (for  $\tau < 1$ ) and  $\delta_{MLE}$  are strictly smaller than  $\sigma_1^2$  for all  $\Delta \geq 0$  with their limits, as  $\Delta \rightarrow \infty$ , equal to  $\sigma_1^2$ . For  $\delta_{MIN}$  with  $\tau = 1$ , the MSE equals  $\sigma_1^2$  for  $\Delta = 0$ . For  $\delta_P$  equality holds for  $\Delta = 0$  as well as for  $\Delta \rightarrow \infty$ .

The previous theorem proves the conjectures 2 and 4 (for MIN only when  $\tau \leq 1$ ). The next theorem shows that the WLE, the MLE as well as MIN are inadmissible and a class of dominators of the MLE is given. (Such dominators can be found in Shao and Strawderman (1996)). Thus the next result answers Question 1, 2 and 3, the latter for MIN.

**Theorem 4.3** The estimators  $\delta_{WLE}$ ,  $\delta_{MLE}$  and  $\delta_{MIN}$  are inadmissible. Further,  $\delta_{MLE}$  is dominated by

$$\frac{\tau Y_1 + Y_2}{1+\tau} - \delta_2^* \left( \frac{Y_2 - Y_1}{1+\tau} \right)$$

where  $\delta_2^*$  is a dominator of the maximum likelihood estimator of a non-negative normal mean based on a single observation.



In the following theorem we state a result of Cohen and Sackrowitz (1970) concerning the admissibility and minimaxity of PIT. In the Appendix we give a simpler proof of the admissibility. A simpler proof of the minimaxity of PIT can be found in Kumar and Sharma (1988, Theorem 2.3). That theorem thus proves Conjecture 4 for PIT and it answers Question 3 affirmatively for that estimator.

**Theorem 4.4** *The Pitman estimator is admissible and minimax. The minimax value for our problem equals  $\sigma_1^2$ .*

The following theorem contains more minimaxity results and proves Conjecture 1 for the MLE as well as Conjecture 4 for the MIN when  $\tau \leq 1$ .

**Theorem 4.5** *The estimators  $\delta_{MLE}$  and  $\delta_{WLE}$  are minimax and so is  $\delta_{MIN}$  when  $\tau \leq 1$ . Further,  $\delta_{MIN}$  is not minimax when  $\tau > 1$ .*

## 5 Discussion

In this article we have tried to show how the intuitively natural idea of the relevance weighted likelihood can be used in parametric estimation to trade bias for precision and thereby reduce the MSE in fortuitous circumstances. The resulting estimators use all the relevant information and not just the direct sample information from the population of interest. By comparing those estimators with others that were obtained earlier for the same purpose we find the weighted likelihood to be promising.

Although we demonstrate the value of our method in a specific normal means estimation context the method itself has wide applicability. Methods of the type described here seem likely to assume increasingly greater importance as the space-time scales of modern experiments continue to expand thanks to need and technological feasibility. Indeed the classical repeated sampling paradigm on which Fisher bases his theory of the MLE will become increasingly untenable as that scale grows. Reliance on biased but relevant sample data will become increasingly imperative.

To conclude we summarize the results of our investigation in Section 4 of the conjectures and questions suggested by the numerical work in Section 3.

**Question's 1 and 2.** We answer negatively these questions on the admissibility of the MLE and WLE respectively (in Theorem 4.3).

- Question 3.** Theorem 4.3 gives a negative answer on the admissibility of MIN when  $\sigma_1 = \sigma_2$ . However, Theorem 4.4 answers positively the same question for PIT.
- Question 4.** This question on the form of the MIN when the population variances are unequal remains open.
- Conjecture 1.** Theorem 4.5 proves the claim in this conjecture that the ULE and MLE are minimax.
- Conjecture 2.** Theorem 4.2 proves the claim here that the ULE is inadmissible and dominated by the MLE.
- Conjecture 3.** The truth of the claim in the conjecture has not been established; we do not now whether or not the WLE is minimax in spite of the very strong numerical evidence against it.
- Conjecture 4.** Theorem 4.2 proves both parts of this conjecture at least that the MIN dominates the ULE when  $r \leq 1$  and the PIT in any case.
- Conjecture 5.** Theorem 4.5 proves the statement in this conjecture for the MIN. We show it to be minimax when  $\sigma_2 \leq \sigma_1$ . At the same time, Theorem 4.4 shows that PIT is minimax whatever be the  $\sigma$ 's.

## BIBLIOGRAPHY

- Akaike, H. (1973). Information theory and an extension of entropy maximization principle. In *2nd International Symposium on Information Theory*, eds. B. N. Petrov and F. Csak, Kiado: Akademia, p. 276-281.
- Akaike, H. (1977). On entropy maximization principle. In *Applications of Statistics*, ed. P.R. Krishnalah Amsterdam: Noeth-Holland, 27-41.
- Akaike, H. (1978). A Bayesian analysis if the minimum AIC procedure. *Ann. Inst. Statist. Math.*, 30A, 9-14.
- Akaike, H. (1982). On the fallacy of the likelihood principle. *Statist. Probab. Lett.* 1, 75-78.
- Akaike, H. (1983). Information measures and model selection. *Proc. 44th ISI Session*, 1, 277-291.
- Akaike, H. (1985). Prediction and entropy. *A Celebration of Statistics, The ISI Centenary Volume*. Berlin, Spring-Verlag.

- Al-Saleh, M. F. (1997). Estimating the mean of a normal population utilizing some available information: a bayesian approach. *J. Inform. Optim. Sci.*, 18, 1-7.
- Berger, J.O. (1985). *Statistical Decision Theory and Bayesian Analysis*, Second edition. New York: Springer - Verlag.
- Blumenthal, S. and Cohen, A. (1968a). Estimation of the larger translation parameter. *Ann. Math. Statist.*, 39, 502-516.
- Blumenthal, S. and Cohen, A. (1968b). Estimation of two ordered translation parameters. *Ann. Math. Statist.*, 39, 517-530.
- Blyth, C.R. (1951). On minimax statistical decision procedures and their admissibility. *Ann. Math. Statist.*, 22, 22-42.
- Brown, L.D. (1986). *Fundamentals of Statistical Exponential Families*. Lecture Notes-Monograph Series, Vol. 9, Institute of Mathematical Statistics, Hayward, California.
- Cohen, A. and Sackrowitz, H. B. (1970). Estimation of the last mean of a monotone sequence. *Ann. Math. Statist.*, 41, 2021-2034.
- Efron, B. (1996). Empirical Bayes methods for combining likelihoods (with discussion). *J. Amer. Statist. Assoc.*, 91, 538-565.
- Fan, J., Heckman, N.E. and Wand, W.P. (1995). Local polynomial kernel regression for generalized linear models and quasi-likelihood functions. *J. Amer. Statist. Assoc.*, 90, 141-150.
- Hu, F. (1994). Relevance weighted smoothing and a new bootstrap method. Ph.D. Thesis, Department of Statistics, University of British Columbia.
- Hu, F. (1997). Asymptotic properties of relevance weighted likelihood estimations. *Canad. J. Statist.*, 25, 45-60.
- Hu, F. and Zidek, J.V. (1997). The relevance weighted likelihood. *Revised to J. Amer. Statist. Assoc.*
- Katz, M.W. (1961). Admissible and minimax estimates of parameters in truncated spaces. *Ann. Math. Statist.*, 32, 136-142.
- Kumar, S. and Sharma, D. (1988). Simultaneous estimation of ordered parameters. *Comm. Statist. Theory Methods*, 17, 4315-4336.
- Kumar, S. and Sharma, D. (1993). Minimavity of the Pitman estimator of ordered normal means when the variances are unequal. *J. Indian Soc. Agricultural Statist.*, 45, 230-234.



Lee, C-I. C. (1981). The quadratic loss of isotonic regression under normality. *Ann. Statist.*, 9, 686-688.

Shao, P.Y-S. and Strawderman, W.E. (1996). Improving on the MLE of a positive normal mean. *Statist. Sinica*, 6, 259-274.

Staniswalis, J.G. (1989). The kernel estimate of a regression function in likelihood-based models. *J. Amer. Statist. Assoc.*, 84, 276 - 283.

Stigler, S M (1990). The 1988 Neyman Memorial Lecture: A Galtonian perspective on shrinkage estimators. *Statist. Sci.*, 5, 147-155.

Tibshirani, R. and Hastie, T. (1987). Local likelihood estimation. *J. Amer. Statist. Assoc.* 82, 559-567.

van Eeden, C. (1995). Minimax estimation of a lower-bounded scale parameter of a gamma distribution for scale-invariant squared-error loss. *Canad. J. Statist.*, 23, 245-256.

## A Appendix

This Appendix contains the proofs of the results presented in Sections 4. In these proofs the following four results (stated in the form of lemmas and a corollary) are used.

The first lemma contains the well-known Stein identity.

**Lemma A.1** For a  $\mathcal{N}(\nu, \gamma^2)$  random variable  $Z$  and a function  $g$  which is almost everywhere (with respect to Lebesgue measure) differentiable

$$\mathcal{E}(Z - \nu)g(Z) = \gamma^2 \mathcal{E}g'(Z). \quad (\text{A.1})$$

The following corollary gives expressions for the mean-squared-error of the estimator  $Y_1 + \varphi(W)$  of  $\theta_1$ . These expressions follow immediately from Lemma A.1 and the fact that the distribution of  $Y_1$ , conditional on  $W$ , is given by

$$\mathcal{N}\left(\theta_1 + \frac{\Delta - W}{1 + \tau}, \frac{\sigma_2^2}{1 + \tau}\right) = \mathcal{N}\left(\frac{\sigma_2^2 \theta_1 + \sigma_1^2 (\theta_2 - W)}{\sigma_1^2 + \sigma_2^2}, \frac{\sigma_1^2 \sigma_2^2}{\sigma_1^2 + \sigma_2^2}\right). \quad (\text{A.2})$$

**Corollary A.1** The mean-squared-error of the estimator  $Y_1 + \varphi(W)$  of  $\theta_1$  is given by

$$\mathcal{E}_\theta(Y_1 - \theta_1 + \varphi(W))^2 = \sigma_1^2 + 2\mathcal{E}_\theta(Y_1 - \theta_1)\varphi(W) + \mathcal{E}_\theta\varphi^2(W),$$



where

$$\begin{aligned}\mathcal{E}_\theta(Y_1 - \theta_1)\varphi(W) &= \mathcal{E}_\theta\{\varphi(W)(\mathcal{E}_\theta(Y_1 - \theta_1) | W)\} \\ &= \frac{1}{1+\tau} \mathcal{E}_\theta(\Delta - W)\varphi(W).\end{aligned}\tag{A.3}$$

Further, if  $\varphi(W)$  is differentiable almost everywhere,

$$\mathcal{E}_\theta(Y_1 - \theta_1)\varphi(W) = -\sigma^2 \mathcal{E}_\theta\varphi'(W).\tag{A.4}$$

In our next lemma, a rotation technique used by Blumenthal and Cohen (1968a) (see also Cohen and Sackrowitz (1970)) is applied.

**Lemma A.2** *Let*

$$\begin{aligned}X_1 &= \frac{\tau Y_1 + Y_2}{1+\tau} & X_2 &= \frac{-Y_1 + Y_2}{1+\tau} \\ \mu_1 &= \mathcal{E}_\theta X_1 = \frac{\tau\theta_1 + \theta_2}{1+\tau} & \mu_2 &= \mathcal{E}_\theta X_2 = \frac{-\theta_1 + \theta_2}{1+\tau}\end{aligned}\tag{A.5}$$

Then  $Y_1 + \varphi(W)$  is inadmissible for estimating  $\theta_1$  based on  $(Y_1, Y_2)$  under the condition  $\theta_1 \leq \theta_2$  if  $\delta_2(X_2) = X_2 - \varphi((1+\tau)X_2)$  is inadmissible for estimating  $\mu_2$  based on  $X_2$  under the condition  $\mu_2 \geq 0$ . Further, if  $\delta_2^*(X_2)$  dominates  $\delta_2(X_2)$  for estimating  $\mu_2$  under the condition  $\mu_2 \geq 0$  based on  $X_2$ , then  $X_1 - \delta_2^*(X_2)$  dominates  $Y_1 + \varphi(W)$  for estimating  $\theta_1$  under the condition  $\theta_1 \leq \theta_2$  based on  $(Y_1, Y_2)$ .

*Proof.* First note that, under  $\theta_1 \leq \theta_2$ ,  $\mu_1$  is unrestricted while  $\mu_2 \geq 0$ . Further

$$Y_1 + \varphi(W) = X_1 - (X_2 - \varphi((1+\tau)X_2)) = X_1 - \delta_2(X_2).$$

The inadmissibility of  $\delta_2(X_2)$  for estimating  $\mu_2 \geq 0$  based on  $X_2$  implies that there exists an estimator  $\delta_2^*(X_2)$  which dominates  $\delta_2(X_2)$  on  $\{\mu_2 | \mu_2 \geq 0\}$ . But this is easily shown to imply that

$$X_1 - \delta_2^*(X_2) = \frac{\tau Y_1 + Y_2}{1+\tau} - \delta_2^*\left(\frac{Y_2 - Y_1}{1+\tau}\right)$$

dominates

$$X_1 - \delta_2(X_2) = \frac{\tau Y_1 + Y_2}{1+\tau} - \delta_2\left(\frac{Y_2 - Y_1}{1+\tau}\right) = Y_1 + \varphi(W)$$

as an estimator of  $\theta_1$  under the condition  $\theta_1 \leq \theta_2$ .  $\square$

The following result is used several times in our proofs.

**Lemma A.3**

$$\mathcal{E}_\theta WI(W < 0) = \Delta \left(1 - \Phi\left(\frac{\Delta}{\sigma}\right)\right) - \sigma \phi\left(\frac{\Delta}{\sigma}\right),$$

$$\mathcal{E}_\theta W^2 I(W < 0) = (\Delta^2 + \sigma^2) \left(1 - \Phi\left(\frac{\Delta}{\sigma}\right)\right) - \Delta \sigma \phi\left(\frac{\Delta}{\sigma}\right).$$

*Proof.* We have

$$\begin{aligned} \mathcal{E}_\theta WI(W < 0) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^0 w e^{-\frac{(w-\Delta)^2}{2\sigma^2}} dw \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{-\Delta/\sigma} (\Delta + y\sigma) e^{-\frac{y^2}{2}} dy \\ &= \Delta \left(1 - \Phi\left(\frac{\Delta}{\sigma}\right)\right) - \sigma \phi\left(\frac{\Delta}{\sigma}\right). \end{aligned}$$

Further,

$$\begin{aligned} \mathcal{E}_\theta W^2 I(W < 0) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\Delta/\sigma} (\Delta + y\sigma)^2 e^{-\frac{y^2}{2}} dy = \\ &\Delta^2 \Phi\left(-\frac{\Delta}{\sigma}\right) - 2\sigma \Delta \phi\left(-\frac{\Delta}{\sigma}\right) \\ &\quad - \sigma^2 \left(-\frac{\Delta}{\sigma} \phi\left(-\frac{\Delta}{\sigma}\right) - \Phi\left(-\frac{\Delta}{\sigma}\right)\right) = \\ &(\Delta^2 + \sigma^2) \left(1 - \Phi\left(\frac{\Delta}{\sigma}\right)\right) - \Delta \sigma \phi\left(\frac{\Delta}{\sigma}\right). \end{aligned}$$

□

We are now ready to give the proofs of the results in Section 4.

PROOF OF THE FORMULA FOR THE PITMAN ESTIMATOR  $\delta_P$  (see (3.5))

The Pitman estimator  $\delta_P(Y_1, Y_2)$  is given by

$$\delta_P(Y_1, Y_2) = Y_1 + \frac{I_1}{I_2}$$

where

$$\begin{aligned} 2\pi\sigma_1\sigma_2 I_1 &= \\ \int \int_{\theta_1 \leq \theta_2} (\theta_1 - Y_1) \exp\left(-\frac{(\theta_1 - Y_1)^2}{2\sigma_1^2}\right) \exp\left(-\frac{(\theta_2 - Y_2)^2}{2\sigma_2^2}\right) d\theta_1 d\theta_2 \end{aligned}$$

and

$$2\pi\sigma_1\sigma_2 I_2 = \int \int_{\theta_1 \leq \theta_2} \exp\left(-\frac{(\theta_1 - Y_1)^2}{2\sigma_1^2}\right) \exp\left(-\frac{(\theta_2 - Y_2)^2}{2\sigma_2^2}\right) d\theta_1 d\theta_2.$$

Note that  $I_2 = P(\theta_1 \leq \theta_2)$  with  $\theta_1$  and  $\theta_2$  independent and  $\theta_1 \sim \mathcal{N}(Y_1, \sigma_1^2)$ ,  $\theta_2 \sim \mathcal{N}(Y_2, \sigma_2^2)$ . So

$$I_2 = P\left(\frac{\theta_1 - Y_1 - (\theta_2 - Y_2)}{\sigma} \leq \frac{W}{\sigma}\right) = \Phi\left(\frac{W}{\sigma}\right), \quad (\text{A.6})$$

where  $\sigma^2 = \sigma_1^2 + \sigma_2^2$ .

Further,  $I_1$  can be represented through the expression

$$\begin{aligned} 2\pi\sigma_1\sigma_2 I_1 &= \\ -\sigma_1^2 \int_{-\infty}^{\infty} \exp\left(-\frac{(\theta_2 - Y_2)^2}{2\sigma_2^2}\right) \int_{-\infty}^{\theta_2} \frac{d}{d\theta_1} \exp\left(-\frac{(\theta_1 - Y_1)^2}{2\sigma_1^2}\right) d\theta_1 d\theta_2 &= \\ -\sigma_1^2 \int_{-\infty}^{\infty} \exp\left(-\frac{(\theta_2 - Y_2)^2}{2\sigma_2^2} - \frac{(\theta_2 - Y_1)^2}{2\sigma_1^2}\right) d\theta_2 &= \\ -\sigma_1^2 \sigma_2 \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}x^2 - \frac{(x\sigma_2 + W)^2}{2\sigma_1^2}\right) dx, \end{aligned}$$

where

$$x^2 + \frac{(x\sigma_2 + W)^2}{\sigma_1^2} = \frac{\sigma^2}{\sigma_1^2} \left(z + \frac{\sigma_2 W}{\sigma^2}\right)^2 + \frac{W^2}{\sigma^2}.$$

So,

$$\begin{aligned} I_1 &= \\ \frac{-\sigma_1}{2\pi} \exp\left(-\frac{1}{2}\left(\frac{W}{\sigma}\right)^2\right) \int_{-\infty}^{\infty} \exp\left(-\frac{\sigma^2}{2\sigma_1^2} \left(z + \frac{\sigma_2 W}{\sigma^2}\right)^2\right) dz &= \\ -\sigma_1 \phi\left(\frac{W}{\sigma}\right) \frac{\sigma_1}{\sigma} &= -\frac{\sigma_1^2}{\sigma} \phi\left(\frac{W}{\sigma}\right). \end{aligned} \quad (\text{A.7})$$

The result follows immediately from the above expressions for  $I_1$  and  $I_2$ .

#### PROOF OF THEOREM 4.1

1. The two formulas for the MSE of  $\delta_{WLE}$  can be obtained by using Corollary A.1 with (A.3) and (A.4) respectively for  $\mathcal{L}_\theta(Y_1 - \theta_1)\varphi(W) = \mathcal{L}_\theta W \dot{\alpha}(W)$ ;

2. The first expression for the MSE of  $\delta_{MLE}$  is obtained from Corollary A.1 using (A.3) for  $\mathcal{E}_\theta(Y_1 - \theta_1) \varphi(W) = (1 + \tau)^{-1} \mathcal{E}_\theta(Y_1 - \theta_1) W$ . The second expression is obtained from Lemma A.3;
3. The first expression for the MSE of  $\delta_{MIN}$  follows from Corollary A.1 using (A.3) for  $\mathcal{E}_\theta(Y_1 - \theta_1) \varphi(W) = \mathcal{E}_\theta(Y_1 - \theta_1) W$ . The second expression is obtained from Lemma A.3;
4. For the MSE of the Pitman estimator, note that

$$\varphi(W) = -\frac{\sigma_1^2}{\sigma \sigma_2} \frac{\phi(W'')}{\Phi(W'')},$$

so

$$\mathcal{E}_\theta(Y_1 - \theta_1) \varphi(W) = -\frac{\sigma_1^2}{\sigma \sigma_2} \mathcal{E}_\theta(Y_1 - \theta_1) \frac{\phi(W'')}{\Phi(W'')},$$

where  $W'' = W/\sigma$ . Further, using Stein's identity (see Lemma A.1) and the fact that

$$\frac{d}{dx} \frac{\phi(x)}{\Phi(x)} = -\left(x \frac{\phi(x)}{\Phi(x)} + \left(\frac{\phi(x)}{\Phi(x)}\right)^2\right), \quad (\text{A.8})$$

gives

$$\begin{aligned} \mathcal{E}_\theta(Y_1 - \theta_1) \frac{\phi(W'')}{\Phi(W'')} &= \mathcal{E}_\theta \left[ \mathcal{E}_\theta \left\{ (Y_1 - \theta_1) \frac{\phi(Y_2'' - Y_1'')}{\Phi(Y_2'' - Y_1'')} \mid Y_2 \right\} \right] = \\ \sigma_1^2 \mathcal{E}_\theta \left[ \mathcal{E}_\theta \left\{ \frac{d}{dY_1} \frac{\phi(Y_2'' - Y_1'')}{\Phi(Y_2'' - Y_1'')} \mid Y_2 \right\} \right] &= \\ \frac{\sigma_1^2}{\sigma} \mathcal{E}_\theta \left\{ (Y_2'' - Y_1'') \frac{\phi(Y_2'' - Y_1'')}{\Phi(Y_2'' - Y_1'')} + \left( \frac{\phi(Y_2'' - Y_1'')}{\Phi(Y_2'' - Y_1'')} \right)^2 \right\}, \end{aligned}$$

where, for  $i = 1, 2$ ,  $Y_i'' = Y_i/\sigma$ . So,

$$\begin{aligned} MSE - \sigma_1^2 &= \\ 2\mathcal{E}_\theta(Y_1 - \theta_1) \varphi(W) + \mathcal{E}_\theta \varphi^2(W) &= \\ -2 \frac{\sigma_1^2}{\sigma} \left[ \frac{\sigma_1^2}{\sigma} \mathcal{E}_\theta \left\{ W'' \frac{\phi(W'')}{\Phi(W'')} + \left( \frac{\phi(W'')}{\Phi(W'')} \right)^2 \right\} \right] & \quad (\text{A.9}) \\ + \frac{\sigma_1^4}{\sigma^2} \mathcal{E}_\theta \left( \frac{\phi(W'')}{\Phi(W'')} \right)^2 &= \end{aligned}$$



$$-\frac{\sigma_1^4}{\sigma^2} \left\{ \mathcal{E}_\theta \left( \frac{\phi(W'')}{\Phi(W'')} \right)^2 + 2\mathcal{E}_\theta W'' \frac{\phi(W'')}{\Phi(W'')} \right\}.$$

Now use

$$\mathcal{E}_\theta W'' \frac{\phi(W'')}{\Phi(W'')} = \mathcal{E}_\theta \left( \frac{W - \Delta}{\sigma} \right) \frac{\phi(W'')}{\Phi(W'')} + \frac{\Delta}{\sigma} \mathcal{E}_\theta \frac{\phi(W'')}{\Phi(W'')}$$

and apply the Stein identity to the right hand side. This gives, using (A.8),

$$\mathcal{E}_\theta W'' \frac{\phi(W'')}{\Phi(W'')} = -\mathcal{E}_\theta W'' \frac{\phi(W'')}{\Phi(W'')} - \mathcal{E}_\theta \left( \frac{\phi(W'')}{\Phi(W'')} \right)^2 + \frac{\Delta}{\sigma} \mathcal{E}_\theta \frac{\phi(W'')}{\Phi(W'')}$$

or

$$2\mathcal{E}_\theta W'' \frac{\phi(W'')}{\Phi(W'')} = -\mathcal{E}_\theta \left( \frac{\phi(W'')}{\Phi(W'')} \right)^2 + \frac{\Delta}{\sigma} \mathcal{E}_\theta \frac{\phi(W'')}{\Phi(W'')}. \quad (\text{A.10})$$

Then using (A.9) and (A.10) gives

$$MSE = \sigma_1^2 - \frac{\sigma_1^4}{\sigma^2} \Delta \mathcal{E}_\theta \frac{\phi(W'')}{\Phi(W'')}.$$

The above proof is a generalization of Al-Saleh's (1997) proof for the case where  $\sigma_1 = \sigma_2$ . Another, similar, proof for the case where  $\tau = 1$  can be obtained from Kumar and Sharma (1993). They use a weighted (by the reciprocals of the known variances) squared-error loss function for estimating  $(\theta_1, \theta_2)$  when  $\theta_1 \leq \theta_2$  (In fact, they consider the more general problem of estimating a  $k$ -dimensional parameter under a complete order restriction). From their formula for the MSE of their Pitman estimator, ours (and thus Al-Saleh's) formula for the MSE of  $\delta_P$  can be obtained by using (A.6) and (A.7).

#### PROOF OF THEOREM 4.2

That  $\delta_{MLE}$  dominates  $\delta_{VLE}$  follows from a result of Lee (1981). He shows that for independent  $Y_i \sim \mathcal{N}(\theta_i, 1)$ ,  $i = 1, \dots, k$ , with  $\theta_1 \leq \dots \leq \theta_k$ , the  $i$ -th component of the order-restricted MLE dominates  $Y_i$ ,  $i = 1, \dots, k$ . For our particular case, where  $k = 2$ , the result can more easily be proved by using the second line of (4.2) and the following inequalities

$$\Delta \mathcal{E}_\theta W I(W < 0) \leq 0 \text{ for all } \Delta \geq 0 \quad (\text{A.11})$$

$$\mathcal{E}_\theta W^2 I(W < 0) > 0 \text{ for all } \Delta \geq 0. \quad (\text{A.12})$$

From (4.5) it is immediately clear that  $\delta_P$  dominates  $\delta_{VLE}$ .

To see that  $\delta_{MIN}$  dominates  $\delta_{WLE}$  when  $\tau \leq 1$ , note that (see the second line of (4.4))

$$\Delta \mathcal{E}_\theta W I(W < 0) \leq 0 \text{ for all } \Delta \geq 0 \quad (\text{A.13})$$

$$(1 - \tau) \mathcal{E}_\theta W^2 I(W < 0) = 0 \text{ for all } \Delta \geq 0 \text{ when } \tau = 1 \quad (\text{A.14})$$

$$(1 - \tau) \mathcal{E}_\theta W^2 I(W < 0) > 0 \text{ for all } \Delta \geq 0 \text{ when } \tau < 1. \quad (\text{A.15})$$

From the second line of (4.4) and the inequalities (A.13) and (A.15) it follows that the MSE of  $\delta_{MIN}$  with  $\tau < 1$  is strictly smaller than  $\sigma_1^2$  for all  $\Delta \geq 0$ . For  $\delta_{MLE}$ , the second line of (4.2) and the inequalities (A.11) and (A.12) imply that its MSE is strictly smaller than  $\sigma_1^2$  for all  $\Delta \geq 0$ . As for the limits, as  $\Delta \rightarrow \infty$ , of the MSEs of  $\delta_{MIN}$  and  $\delta_{MLE}$ , use the last line of (4.4) and of (4.2) and note that

$$\lim_{\Delta \rightarrow \infty} \Delta^2 \left( 1 - \Phi \left( \frac{\Delta}{\sigma} \right) \right) = 0$$

and

$$\lim_{\Delta \rightarrow \infty} \Delta \phi \left( \frac{\Delta}{\sigma} \right) = 0.$$

That, for  $\tau = 1$ , the MSE of  $\delta_{MIN}$  equals  $\sigma_1^2$  for  $\Delta = 0$  follows immediately from the last line of (4.4).

Finally (using (4.5)), the MSE of  $\delta_\rho$  clearly equals  $\sigma_1^2$  when  $\Delta = 0$ . That its MSE converges to  $\sigma_1^2$  when  $\Delta \rightarrow \infty$  can be seen from (4.5) by noting that

$$\Delta \mathcal{E}_\theta \frac{\phi \left( \frac{W}{\sigma} \right)}{\Phi \left( \frac{W}{\sigma} \right)} = \Delta \mathcal{E} \frac{\phi \left( Z + \frac{\Delta}{\sigma} \right)}{\Phi \left( Z + \frac{\Delta}{\sigma} \right)},$$

where  $Z \sim \mathcal{N}(0, 1)$ . The result then follows from the fact that, for each fixed  $z$ ,  $\Delta \phi(z + \Delta/\sigma)/\Phi(z + \Delta/\sigma)$  is bounded in  $z$  for  $\Delta \geq 0$  and converges to zero as  $\Delta \rightarrow \infty$ .

#### PROOF OF THEOREM 4.3

That  $\delta_{WLE}$  is inadmissible can be shown by using Lemma A.2 with (see (3.1))  $\phi(W) = W \alpha(W)$ . Then

$$\delta_2(X_2) = X_2 - \phi((1 + \tau)X_2) = X_2 - X_2 \sigma^2 \left( \sigma^2 + (1 + \tau)^2 (\max(0, X_2))^2 \right)^{-1}.$$

This estimator does not satisfy Brown's (1986, Theorem 4.23) necessary conditions for admissibility for estimating  $\mu_2 \geq 0$  based on  $X_2$ .

For the proof of the inadmissibility of  $\delta_{MLE}$  again use Lemma A.2, this time with  $\varphi(W) = (1 + \tau)^{-1} W_-$ . Then  $\delta_2(X_2) = X_2 - \varphi((1 + \tau)X_2) = \max(0, X_2)$ . This  $\delta_2(X_2)$

is the maximum likelihood estimator for estimating  $\mu_2 \geq 0$  based on  $X_2$  which is well-known to be inadmissible for squared loss. The formula for the dominators of  $\delta_{MLE}$  follows directly from Lemma A.2.

For the proof of the inadmissibility of  $\delta_{MLN}$ , use Lemma A.2 with  $\delta_2(X_2) = X_2 I(X_2 \geq 0) - \tau X_2 I(X_2 < 0)$ . That this  $\delta_2$  is inadmissible for estimating  $\mu_2 \geq 0$  based on  $X_2$  follows from Theorem 4.23 of Brown (1986) and the fact that  $\delta_2(X_2)$  is not monotone in  $X_2$ .

#### PROOF OF THEOREM 4.4

As already noted above, Kumar and Sharma (1988, Theorem 2.3) give a proof of the minimaxity of  $\delta_P$ . Their proof is very much simpler than the one given by Cohen and Sackrowitz (1970). The Kumar-Sharma proof is based on an extension of a result of Blumenthal and Cohen (1968b, Theorem 3.0).

For an alternate and simpler proof of the admissibility of  $\delta_P$ , use the transformation (A.5). Then (see the proof of Lemma A.2)

$$\delta_P(Y_1, Y_2) = X_1 - \delta_2(X_2),$$

where

$$\delta_2(X_2) = X_2 + \sigma(X_2) \frac{\phi\left(\frac{X_2}{\sigma(X_2)}\right)}{\Phi\left(\frac{X_2}{\sigma(X_2)}\right)},$$

where  $\sigma^2(X_2)$  is the variance of  $X_2$ . Further,  $\theta_1 = \mu_1 - \mu_2$  and

$$\theta_1 \leq \theta_2 \iff \mu_1 \in (-\infty, \infty), \mu_2 \geq 0.$$

So, it is now sufficient to show that  $\delta(X) = X_1 - \delta(X_2)$  is admissible for estimating  $\mu_1 - \mu_2$  based on  $X = (X_1, X_2)$  when  $\mu_2 \geq 0$ . We will show this by using Blyth's (1951) method.

Suppose that there exists an estimator  $\delta'(X)$  which dominates  $\delta(X)$  on  $\Omega = \{\mu \mid \mu_1 \in (-\infty, \infty), \mu_2 \geq 0\}$ . Then, because the risk function  $R(\delta_\mu, \mu)$  of every estimator  $\delta_\mu(X)$  is continuous in  $\mu$  for  $\mu \in \Omega$ , there exists an  $\varepsilon > 0$  and a rectangle  $S = (\mu_{1,1}, \mu_{1,2}) \times (\mu_{2,1}, \mu_{2,2}) \subset \Omega$  such that

$$R(\delta, \mu) - R(\delta', \mu) > \varepsilon \quad \text{on } S. \quad (\text{A.16})$$

Now take a sequence of priors  $\lambda_n, n = 1, 2, \dots$  for  $\mu \in \Omega$  where, for each  $n$ ,  $\mu_1$  and  $\mu_2$  are independent,  $\mu_1$  with the improper uniform prior on  $(-\infty, \infty)$  and  $\mu_2$  with

density  $e^{\mu_2/n}/n$  on  $\mu_2 \geq 0$ . Then (A.16) implies that

$$r_n(\delta) - r_n(\delta') > \varepsilon(\mu_{1,2} - \mu_{1,1}) \left( e^{-\frac{\varepsilon_{1,1}}{\varepsilon}} - e^{-\frac{\varepsilon_{1,2}}{\varepsilon}} \right) = O\left(\frac{1}{n}\right), \quad (\text{A.17})$$

where the  $r_n$ 's are Bayes risks with respect to  $\lambda_n$ .

Now let, for  $i = 1, 2$ ,  $\delta_{n,i}(X_i)$  be the Bayes estimator of  $\mu_i$  based on  $X_i$  with respect to the (marginal) prior of  $\mu_i$ . Then, by the prior independence of  $\mu_1$  and  $\mu_2$  and the conditional independence of  $X_1$  and  $X_2$  given  $\mu_1$  and  $\mu_2$ , the Bayes estimator of  $\mu_1 - \mu_2$  for the prior  $\lambda_n$  based on  $X$ , is given by

$$\delta_n(X) = \delta_{n,1}(X_1) - \delta_{n,2}(X_2),$$

where (see Katz (1961))

$$\delta_{n,2}(X_2) = X_2 - \frac{1}{n} + \sigma(X_2) \frac{\phi\left(\frac{X_2 - \frac{1}{n}}{\sigma(X_2)}\right)}{\Phi\left(\frac{X_2 - \frac{1}{n}}{\sigma(X_2)}\right)},$$

Further (see Katz, (1961))

$$r_n(\delta) - r_n(\delta_n) = r_{n,2}(\delta_2) - r_{n,2}(\delta_{n,2}) = o\left(\frac{1}{n}\right), \quad (\text{A.18})$$

where  $r_{n,2}$  is the Bayes risk of an estimator based on  $X_2$  with respect to the (marginal) prior of  $\mu_2$ .

But (A.17) and (A.18) imply that, for sufficiently large  $n$ ,

$$\frac{r_n(\delta) - r_n(\delta')}{r_n(\delta) - r_n(\delta_n)} > 1$$

which contradicts the fact that, for each  $n$ ,  $\delta_n(X)$  is the Bayes estimator of  $\mu_2$  with respect to  $\lambda_n$  based on  $X_2$ .

**Remark A.1** Katz's minimaxity and admissibility proofs are incorrect for the general case of the exponential family he considers (see van Eeden (1995)), but the above quoted results of his for the normal mean are correct.

## PROOF OF THEOREM 4.5



The minimaxity results follow immediately from Theorem 4.2 and the fact that, by Theorem 4.4, the minimax value for our problem is equal to the mean square error  $\sigma_1^2$  of  $\delta_{VLE} = Y_1$ . That  $\delta_{MIN}$  is not minimax when  $\tau > 1$  can be seen from the second line of (4.4) by noting that the MSE of  $\delta_{MIN}$  is larger than  $\sigma_1^2$  when  $\Delta = 0$ .