Estimation after Model Selection in a Gaussian Model

Edsel A. Peña

Department of Statistics University of South Carolina E-Mail: pena@stat.sc.edu

Joint with Prof. Vanja Dukić of the Univ. of Chicago.

Research support from the National Science Foundation.

1. Setting and Problem

• Data Given:

$$\boldsymbol{X} \equiv (X_1, X_2, \dots, X_n) \text{ IID } F$$

where F is an unknown distribution function.

• Model \mathcal{M} :

$$F \in \mathcal{M} = \left\{ N(\mu, \sigma^2) : \mu \in \Re, \sigma^2 > 0 \right\}$$

 $N(\mu, \sigma^2) =$ normal distribution with mean μ and variance σ^2 .

• Problems:

- Estimate σ^2 ;
- Given $t \in \Re$, estimate

$$\tau(t) = F(t) = \Pr\{X_1 \le t\} = \Phi\left(\frac{t-\mu}{\sigma}\right).$$

• Well-known (e.g., proved in Stat 201 ... cheers! if I teach it!) that the uniformly minimum variance unbiased (UMVU) estimator of σ^2 is

$$\hat{\sigma}_{UMVU}^2 = S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

• UMVU estimator of $\tau(t)$, derived via Basu's Theorem and

the Lehmann-Scheffe Theorem, is

$$\hat{\tau}_{UMVU}(t) = \mathcal{T}\left(\frac{\sqrt{n-2}z_1(t)}{\sqrt{1-z_1(t)^2}}; n-2\right) I\{|z_1(t)| \le 1\} + I\{z_1(t) > 1\}$$

with

$$z_1(t) = \frac{\sqrt{n}}{n-1} \left(\frac{t-\bar{X}}{S} \right),$$

and $\mathcal{T}(\cdot; k)$ being the Student's *t*-distribution function with k degrees-of-freedom.

• Maximum likelihood (ML) estimator of $\tau(t)$ is

$$\hat{\tau}_{MLE}(t) = \Phi\left(\frac{t-\bar{X}}{S}\right).$$

• Decision-theoretic framework: Loss functions utilized are

$$L_{1}(\hat{\sigma}^{2},(\mu,\sigma^{2})) = \left(\frac{\hat{\sigma}^{2}-\sigma^{2}}{\sigma^{2}}\right)^{2};$$

$$L_{2}(\hat{\tau}(t),(\mu,\sigma^{2})) = (\hat{\tau}(t) - \Phi((t-\mu)/\sigma))^{2};$$

$$L_{3}(\hat{\tau},(\mu,\sigma^{2})) = \int L_{2}(\hat{\tau}(t),(\mu,\sigma^{2}))\Phi((dt-\mu)/\sigma).$$

• Risk functions ("loss functions averaged-out over X")

$$R_1(\hat{\sigma}^2,(\mu,\sigma^2)); \ R_2(\hat{\tau}(t),(\mu,\sigma^2)); \ R_3(\hat{\tau},(\mu,\sigma^2))$$

are the expected values of the respective loss functions with respect to X and when the true parameter values are (μ, σ^2) . Since loss functions are quadratic, then

$$Risk = Variance + Bias^2$$
.

• When using risk functions to evaluate estimators, and if we allow biased estimators, the sample variance S² is *not* the best. It is dominated by the ML and the minimum risk equivariant (MRE) estimators given by:

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \text{ and } \hat{\sigma}_{MRE}^2 = \frac{n}{n+1} \hat{\sigma}_{MLE}^2.$$

- Is the unbiasedness property (i.e., that the average equals the parameter) a 'sacred cow?'
- All nontrivial Bayes estimators are biased ... so you know what will happen if biased estimators are *not* allowed! Also, can sometimes sacrifice some accuracy to gain precision!

Model $\mathcal{M}_0 = \mathcal{M}_{\mu_0}$

• Suppose it is known that $\mu = \mu_0$, so

$$F \in \mathcal{M}_0 = \left\{ N(\mu, \sigma^2) : \mu = \mu_0, \sigma^2 > 0 \right\}.$$

• Under model \mathcal{M}_0 , the appropriate estimators are:

$$\hat{\sigma}_{UMVU}^{2}(\mu_{0}) = \frac{1}{n} \sum_{i=1}^{n} (X_{i} - \mu_{0})^{2};$$

$$\hat{\sigma}_{MRE}^{2}(\mu_{0}) = \frac{1}{n+2} \sum_{i=1}^{n} (X_{i} - \mu_{0})^{2};$$

$$\hat{\tau}_{UMVU}(t;\mu_{0}) = \mathcal{T}\left(\frac{\sqrt{n-1}z_{2}(t)}{\sqrt{1-z_{2}(t)^{2}}};n-1\right) I\{|z_{2}(t)| \leq 1\}$$

$$+I\{z_{2}(t) > 1\};$$

$$z_{2}(t) = \frac{1}{\sqrt{n}} \left(\frac{t-\mu_{0}}{\hat{\sigma}_{UMVU}(\mu_{0})}\right).$$

Estimators developed under *M* are also candidate estimators under *M*₀. Less efficient however since they do *not* exploit the added structure of *M*₀. For instance, under *M*₀,

$$\operatorname{Eff}\left(\hat{\sigma}_{UMVU}^{2}(\mu_{0}):\hat{\sigma}_{UMVU}^{2}\right)=1+\frac{1}{n-1}.$$

Model \mathcal{M}_p : An Intermediate Model

 Suppose instead that we do not know the exact value of μ, but just that it could take one of p possible values. This leads to model M_p:

$$F \in \mathcal{M}_p = \left\{ N(\mu, \sigma^2) : \mu \in \{\mu_1, \dots, \mu_p\}, \sigma^2 > 0 \right\}$$

where $\mu_1, \mu_2, \ldots, \mu_p$ are known constants.

- **Problem:** Under this intermediate model, how should we estimate σ^2 and $\tau(t) = F(t)$? What are the consequences of using the estimators developed under \mathcal{M} , which are also candidate estimators under \mathcal{M}_p ?
- Can we exploit the structure of \mathcal{M}_p to obtain better estimators of σ^2 and $\tau(t)$? What happens when $p \to \infty$?
- Model \mathcal{M}_p can be viewed as having the p sub-models

$$\mathcal{M}_{\mu_1}, \mathcal{M}_{\mu_2}, \ldots, \mathcal{M}_{\mu_p},$$

with σ^2 a common parameter among these sub-models.

2. Motivation and Importance

- Estimation of the variance, the precision parameter, and the distribution function are important from a practical point of view, as well as theoretically.
- Model \$\mathcal{M}_p\$ corresponds to practical settings where there are

 a finite number of possible populations, and a sample \$\mathcal{X}\$
 is taken from one of them. Setting of the Neyman-Pearson
 Lemma and of multiple decision problems.
- C. Stein (1964) developed the approach of hypothesizing that model \mathcal{M}_0 may hold, and by doing *pre-test* to accept or reject this hypothesis and deciding on which estimator to use based on this test, was able to show that $\hat{\sigma}_{MRE}^2$ is also an *inadmissible* estimator of σ^2 !
- Our primary motivating situations leading to this problem are:

- Issue of 'inference after model selection:' What are the consequences of first selecting a sub-model and then performing inference such as estimation or testing hypothesis, with these two steps utilizing the *same* sample data (i.e., *double-dipping*)?
- Survival analysis and reliability settings: Only known that the family of failure times is either Weibull or gamma (also, Cox PH or accelerated failure model). How to estimate the survivor distribution?
- Multiple regression: A subset of predictors is chosen, then other inferences, such as prediction, is performed.
- Smooth goodness-of-fit tests: An embedding approach is used and it is desired to determine the size of the embedding class *adaptively*.
- Others: In nonparametric regression or function estimation, bandwidths in kernel smoothing are determined *adaptively*.

3. Intuitive Strategies

Strategy I: Simply utilize estimators developed under \mathcal{M} , or a fully nonparametric model.

Strategy II (Classical):

Step 1 (Model Selection): Choose the most plausible sub-model using the data $\boldsymbol{X} = (X_1, X_2, \dots, X_n)$.

Step 2 (Inference): Use the best estimators or test procedures in the chosen sub-model, but with these estimators or tests still using the same data X. Resulting procedures become *adaptive*.

Strategy III (Bayesian): Determine adaptively (i.e., using X) the plausibility of each of the sub-models, and form a weighted combination of the sub-model estimators or tests. Resulting procedures are again *adaptive*.

Question: Which strategy leads to better procedures, and how could we justify formally each of these intuitive strategies?

4. Classical Estimators Under \mathcal{M}_p

• Likelihood Function:

$$L(\mu, \sigma^2) = \prod_{i=1}^{p} L_i(\mu_i, \sigma^2)^{I\{\mu=\mu_i\}};$$

For i = 1, 2, ..., p,

$$L_i(\mu_i, \sigma^2) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \left(\frac{1}{\sigma^2}\right)^{n/2} \exp\left\{-\frac{n\hat{\sigma}_i^2}{2\sigma^2}\right\};$$
$$\hat{\sigma}_i^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \mu_i)^2.$$

• Model Selector: $\hat{M} = \hat{M}(X_1, X_2, \dots, X_n)$

$$\hat{M} = \arg \max_{1 \le i \le p} L_i(\mu_i, \hat{\sigma}_i^2) = \arg \min_{1 \le i \le p} \hat{\sigma}_i^2 = \arg \min_{1 \le i \le p} |\bar{X} - \mu_i|.$$

- \hat{M} chooses the sub-model leading to the smallest estimate of σ^2 , or the sub-model whose mean is closest to the sample mean.
- Selector also has the interpretation of being the 'highest posterior probability (hpp)' model selector associated with a noninformative prior on (μ, σ^2) .
- Other selectors could also be utilized! E.g., Akaike's AIC, Schwarz Bayesian information criterion (BIC).

• MLE of (μ, σ^2) under \mathcal{M}_p :

$$(\hat{\mu}_{p,MLE}, \hat{\sigma}_{p,MLE}^2) = (\hat{\mu}_{\hat{M}}, \hat{\sigma}_{\hat{M}}^2) = \sum_{i=1}^p I\{\hat{M} = i\}(\mu_i, \hat{\sigma}_i^2),$$

so the MLE of σ^2 is

$$\hat{\sigma}_{p,MLE}^2 = \hat{\sigma}_{\hat{M}}^2 = \sum_{i=1}^p I\{\hat{M} = i\}\hat{\sigma}_i^2.$$

- A two-stage adaptive estimator of the 'classical' form.
- **Remark:** If $\hat{M} = i$, the adaptive estimator $\hat{\sigma}_{\hat{M}}^2$ does not have the same properties as *i*th sub-model estimator $\hat{\sigma}_i^2$.
- An alternative estimator of the same flavor as above is to use the sub-model's MRE's given by

$$\hat{\sigma}_{MRE,i}^2 = \frac{n}{n+2}\hat{\sigma}_i^2, \quad i = 1, 2, \dots, p,$$

to obtain

$$\hat{\sigma}_{p,MRE}^2 = \hat{\sigma}_{MRE,\hat{M}}^2 = \sum_{i=1}^p I\{\hat{M} = i\}\hat{\sigma}_{MRE,i}^2.$$

 Remark: Label 'MRE' is a misnomer since this estimator need not be the minimum risk equivariant estimator under model *M_p*. • Adaptive estimator (of 'classical form') of $\tau(t) = F(t)$ un-

der model \mathcal{M}_p is:

$$\hat{\tau}_{p,MLE}(t) = \Phi\left(\frac{t-\mu_{\hat{M}}}{\hat{\sigma}_{\hat{M}}}\right) = \sum_{i=1}^{p} I\{\hat{M}=i\}\Phi\left(\frac{t-\mu_{i}}{\hat{\sigma}_{i}}\right)$$

• Another adaptive estimator obtained from the UMVUs of sub-models is as follows: Let

$$z_{3i}(t) = \frac{1}{\sqrt{n}} \left(\frac{t - \mu_i}{\hat{\sigma}_i} \right), \quad i = 1, 2, \dots, p,$$

and for i = 1, 2, ..., p,

$$\hat{\tau}_{UMVU,i}(t) = \mathcal{T}\left(\frac{\sqrt{n-1}z_{3i}(t)}{\sqrt{1-z_{3i}(t)^2}}; n-1\right) I\{|z_{3i}(t)| \le 1\} + I\{z_{3i}(t) > 1\}.$$

An estimator of $\tau(t)$ is

$$\hat{\tau}_{p,UMVU}(t) = \hat{\tau}_{UMVU,\hat{M}}(t) = \sum_{i=1}^{p} I\{\hat{M} = i\}\hat{\tau}_{UMVU,i}(t).$$

• Remark: Notice that the properties of these adaptive estimators are not easily obtainable because of the interplay between the model selector \hat{M} and the sub-model estimator, both of which are using the same sample data. This interplay makes these situations tough to handle.

5. Bayes Estimators Under \mathcal{M}_p

• Joint Prior Distribution for (μ, σ^2) :

$$\pi(\mu, \sigma^2 | \tilde{\boldsymbol{\theta}}, \beta, \kappa) = \left\{ \prod_{i=1}^p \tilde{\theta}_i^{I\{\mu=\mu_i\}} \right\} \frac{\beta^{\kappa-1}}{\Gamma(\kappa-1)} \left(\frac{1}{\sigma^2} \right)^{\kappa} \exp\left(-\frac{\beta}{\sigma^2} \right)$$

- Independent priors between μ and σ^2 , and for σ^2 , an inverted gamma prior.
- Joint Posterior Distribution: By Bayes rule, with $M_{i} = I\{\mu = \mu_{i}\},$ $\pi(\mu, \sigma^{2} \mid \boldsymbol{x}) = C \prod_{i=1}^{p} \left\{ \tilde{\theta}_{i} \left(\frac{1}{\sigma^{2}}\right)^{\frac{n}{2}+\kappa} \exp\left(-\frac{1}{\sigma^{2}} \left[\frac{n\hat{\sigma}_{i}^{2}}{2} + \beta\right]\right) \right\}^{m_{i}};$ $C = \frac{1}{\Gamma(n/2 + \kappa - 1)} \left\{ \sum_{i=1}^{p} \frac{\tilde{\theta}_{i}}{(n\hat{\sigma}^{2}/2 + \beta)^{n/2 + \kappa - 1}} \right\}^{-1}.$
- Posterior Probabilities of Sub-Models:

$$\theta_i(\kappa,\beta,n,\boldsymbol{x}) = \frac{\tilde{\theta}_i(n\hat{\sigma}_i^2/2+\beta)^{-(n/2+\kappa-1)}}{\sum_{j=1}^p \tilde{\theta}_j(n\hat{\sigma}_j^2/2+\beta)^{-(n/2+\kappa-1)}}$$

As n → ∞, and when viewed as a function of X, the posterior probability of the correct sub-model converges almost surely to 1.

• Posterior Density of σ^2 :

$$\pi(\sigma^2 \mid \boldsymbol{x}) = C \sum_{i=1}^{p} \tilde{\theta}_i \left(\frac{1}{\sigma^2}\right)^{-(\kappa+n/2)} \times \exp\left[-\frac{1}{\sigma^2} \left(n\hat{\sigma}_i^2/2 + \beta\right)\right]$$

• Bayes Estimator of σ^2 :

$$\hat{\sigma}_{p,Bayes}^{2}(\kappa,\beta,\tilde{\boldsymbol{\theta}}) = \sum_{i=1}^{p} \theta_{i}(\kappa,\beta,n,\boldsymbol{x}) \times \left\{ \left(\frac{n}{n+2(\kappa-2)}\right) \hat{\sigma}_{i}^{2} + \left(\frac{2(\kappa-2)}{n+2(\kappa-2)}\right) \left(\frac{\beta}{\kappa-2}\right) \right\}.$$

• Estimator is a weighted combination, in contrast to the 'classical forms' of earlier estimators. Note the weights are data-dependent or adaptive!

• Non-Informative Prior:

- Uniform prior for sub-models: $\tilde{\theta}_i = 1/p, i = 1, 2, \dots, p$.
- Jeffrey's prior on each sub-model: $\beta \to 0; \kappa \to 1$.
- The model selector \hat{M} can be interpreted as the highest posterior probability selector corresponding to this noninformative prior distribution.

• Sub-Models Posterior Probabilities:

$$\theta_i^*(n, \boldsymbol{x}) = \frac{(\hat{\sigma}_i^2)^{(-n/2)}}{\sum_{j=1}^p (\hat{\sigma}_j^2)^{(-n/2)}}$$

• Limiting Bayes Estimator of σ^2 :

$$\hat{\sigma}_{p,LB}^2 = \left(\frac{n}{n-2}\right) \sum_{i=1}^p \left\{\frac{(\hat{\sigma}_i^2)^{(-n/2)}}{\sum_{j=1}^p (\hat{\sigma}_j^2)^{(-n/2)}}\right\} \hat{\sigma}_i^2$$

Remark: Estimator actually examined in the sequel did **not** have the multiplier $\left(\frac{n}{n-2}\right)!$

• Sub-Models Limiting Bayes Estimators:

$$\tilde{\sigma}_{LB,i}^2 = \left(\frac{n}{n-2}\right)\hat{\sigma}_i^2, \quad i = 1, 2, \dots, p,$$

• Another adaptive estimator of σ^2 can be formed from these limiting Bayes estimators via

$$\hat{\sigma}_{p,ALB}^2 = \tilde{\sigma}_{LB,\hat{M}}^2 = \left(\frac{n}{n-2}\right)\sum_{i=1}^p I\{\hat{M}=i\}\hat{\sigma}_i^2.$$

Referred to as the *adaptive limiting Bayes estimator*.

• **Remark:** The estimators $\hat{\sigma}_{p,MRE}^2$, $\hat{\sigma}_{p,MLE}^2$, and $\hat{\sigma}_{p,ALB}^2$ belong to the same class of estimators. Consequently, it suffices to derive results for $\hat{\sigma}_{p,MLE}^2$ since results for the other two estimators becomes immediately obtainable.

• Bayes Estimator of $\tau(t)$: Posterior mean of $\tau(t)$. Can

be shown to be

$$\hat{\tau}_{p,Bayes}(t;\kappa,\beta,\boldsymbol{\theta}) = \sum_{i=1}^{p} \theta_i(\kappa,\beta,n,\boldsymbol{x}) \times \mathcal{T}\left(\frac{\sqrt{\kappa-1+\frac{n}{2}}(t-\mu_i)}{\sqrt{\frac{n}{2}\hat{\sigma}_i^2+\beta}}; 2\left(\kappa-1+\frac{n}{2}\right)\right)$$

• For the non-informative prior $(\tilde{\theta}_i = 1/p, \beta \to 0, \kappa \to 1)$, the limiting Bayes estimator of $\tau(t)$ is

$$\hat{\tau}_{p,LB}(t) = \sum_{i=1}^{p} \left\{ \frac{\left(\hat{\sigma}_{i}^{2}\right)^{-n/2}}{\sum_{j=1}^{p} \left(\hat{\sigma}_{j}^{2}\right)^{-n/2}} \right\} \mathcal{T}\left(\frac{t-\mu_{i}}{\hat{\sigma}_{i}}; n\right).$$

- Which is an adaptively weighted combination of the limiting Bayes estimators in the sub-models.
- Bayes framework therefore leads to estimators that are adaptively weighted combinations of the sub-model estimators.
- Classical framework (MLE, for example) produces two-step estimators: Step 1 is the process of choosing the model; and Step 2 is the process of using the estimator in the chosen model.

Recap: Estimators of σ^2 under \mathcal{M}_p

• Developed under \mathcal{M} :

- $\hat{\sigma}_{UMVU}^2 = S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i \bar{X})^2$
- $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i \bar{X})^2$ and $\hat{\sigma}_{MRE}^2 = \frac{n}{n+1} \hat{\sigma}_{MLE}^2$
- Developed under \mathcal{M}_p :
 - $\hat{\sigma}_i^2 = \frac{1}{n} \sum_{j=1}^n (X_j \mu_i)^2, \quad i = 1, 2, \dots, p.$
 - $\hat{M} = \arg\min_{1 \le i \le p} \hat{\sigma}_i^2 = \arg\min_{1 \le i \le p} |\bar{X} \mu_i|$
 - $\hat{\sigma}_{p,MLE}^2 = \hat{\sigma}_{\hat{M}}^2 = \sum_{i=1}^p I\{\hat{M}=i\}\hat{\sigma}_i^2$
 - $\hat{\sigma}_{MRE,i}^2 = \frac{n}{n+2}\hat{\sigma}_i^2, \quad i = 1, 2, \dots, p.$
 - $\hat{\sigma}_{p,MRE}^2 = \hat{\sigma}_{MRE,\hat{M}}^2 = \sum_{i=1}^p I\{\hat{M}=i\}\hat{\sigma}_{MRE,i}^2$
 - $\hat{\sigma}_{p,ALB}^2 = \tilde{\sigma}_{LB,\hat{M}}^2 = \left(\frac{n}{n-2}\right) \sum_{i=1}^p I\{\hat{M}=i\}\hat{\sigma}_i^2$
 - $\hat{\sigma}_{p,LB}^2 = \sum_{j=1}^p \left\{ \frac{(\hat{\sigma}_i^2)^{-(n/2)}}{\sum_{j=1}^p (\hat{\sigma}_j^2)^{-(n/2)}} \right\} \hat{\sigma}_i^2$
- Question: Which among these σ^2 estimators is best in terms of their risk function?

6. Comparison of σ^2 Estimators

- $R\left(\hat{\sigma}_{UMVU}^2, (\mu, \sigma^2)\right) = \frac{2}{n-1}.$
- $R\left(\hat{\sigma}_{MRE}^2, (\mu, \sigma^2)\right) = \frac{2}{n+1}.$
- Efficiency measure relative to $\hat{\sigma}_{UMVU}^2$:

$$\operatorname{Eff}(\hat{\sigma}^2 : \hat{\sigma}^2_{UMVU}) = \frac{R(\hat{\sigma}^2_{UMVU}, (\mu, \sigma^2))}{R(\hat{\sigma}^2, (\mu, \sigma^2))}.$$

• Eff $(\hat{\sigma}_{MRE}^2 : \hat{\sigma}_{UMVU}^2) = \frac{n+1}{n-1} = 1 + \frac{2}{n-1}.$

Properties of \mathcal{M}_p -Based Estimators

• Notation: $Z \sim N(0,1), \mathbf{Z} = (Z_1, Z_2, ..., Z_n)' \sim N_n(\mathbf{0}, \mathbf{I}),$ and $\boldsymbol{\mu} = (\mu_1, \mu_2, ..., \mu_p)'$. With μ_{i_0} the true mean with $i_0 \in \{1, 2, ..., p\}$, let

$$\boldsymbol{\Delta} = \frac{\boldsymbol{\mu} - \mu_{i_0} \boldsymbol{1}}{\sigma}$$

where $\mathbf{1} = (1, 1, \dots, 1)'$.

• **Proposition:** Under \mathcal{M}_p with μ_{i_0} the true mean,

$$\frac{\hat{\sigma}_i^2}{\sigma^2} \stackrel{d}{=} \frac{1}{n} \left(W + V_i^2 \right), i = 1, 2, \dots, p;$$

 $W \sim \chi^2_{n-1}; \boldsymbol{V} = (V_1, \ldots, V_p)' \sim N_p(-\sqrt{n}\boldsymbol{\Delta}, \boldsymbol{J} \equiv \boldsymbol{11}');$

with W and V stochastically independent.

• Representation of **V**:

$$V = Z1 - \sqrt{n}\Delta$$

- Implication: Distributional characteristics of
 ²/σ² depends on (μ, σ²) only through Δ (and n)! This simplifies comparisons, both theoretical and simulated.
- Notation: Given Δ, let Δ₍₁₎ < Δ₍₂₎ < ... < Δ_(p) denote the associated ordered values. Note that Δ always has a zero component.
- **Theorem:** Under \mathcal{M}_p with μ_{i_0} the true mean,

$$\frac{\hat{\sigma}_{p,MLE}^2}{\sigma^2} \stackrel{d}{=} \frac{1}{n} \left\{ W + \sum_{i=1}^p I\{L(\Delta_{(i)}, \boldsymbol{\Delta}) < Z < U(\Delta_{(i)}, \boldsymbol{\Delta})\}(Z - \sqrt{n}\Delta_{(i)})^2 \right\};$$

 $W \sim \chi^2_{n-1}, Z \sim N(0, 1)$, and with $W \amalg Z$; and

$$L(\Delta_{(i)}, \mathbf{\Delta}) = \frac{\sqrt{n}}{2} \left[\Delta_{(i)} + \Delta_{(i-1)} \right];$$

$$U(\Delta_{(i)}, \mathbf{\Delta}) = \frac{\sqrt{n}}{2} \left[\Delta_{(i)} + \Delta_{(i+1)} \right].$$

[Convention: $\Delta_{(0)} = -\infty$ and $\Delta_{(p+1)} = +\infty$.]

• Representation leads to exact mean and variance:

• Mean:

$$EpMLE(\boldsymbol{\Delta}) \equiv E\left\{\frac{\hat{\sigma}_{p,MLE}^2}{\sigma^2}\right\}$$
$$= 1 - \frac{2}{\sqrt{n}} \sum_{i=1}^p \Delta_{(i)} [\phi(L(\Delta_{(i)}, \boldsymbol{\Delta})) - \phi(U(\Delta_{(i)}, \boldsymbol{\Delta}))] + \sum_{i=1}^p \Delta_{(i)}^2 [\Phi(U(\Delta_{(i)}, \boldsymbol{\Delta})) - \Phi(L(\Delta_{(i)}, \boldsymbol{\Delta}))];$$

• Variance:

$$VpMLE(\boldsymbol{\Delta}) \equiv \mathbf{Var} \left\{ \frac{\hat{\sigma}_{p,MLE}^2}{\sigma^2} \right\}$$
$$= \frac{1}{n} \left\{ 2\left(1 - \frac{1}{n}\right) + \frac{1}{n} \left[\sum_{i=1}^p \zeta_{(i)}(4) - \left(\sum_{i=1}^p \zeta_{(i)}(2) \right)^2 \right] \right\};$$

with

$$\xi(k;\Omega_{(i)}) \equiv \mathbf{E}\left\{Z^k I(\Omega_{(i)})\right\} = \int_{L(\Delta_{(i)},\mathbf{\Delta})}^{U(\Delta_{(i)},\mathbf{\Delta})} z^k \phi(z) dz,$$

and, for $m \in \mathbb{Z}_+$,

$$\begin{aligned} \zeta_{(i)}(m) &\equiv \mathbf{E} \left\{ I(\Omega_{(i)}) (Z - \sqrt{n} \Delta_{(i)})^m \right\} \\ &= \sum_{k=0}^m (-1)^{(m-k)} \binom{m}{k} \left(\sqrt{n} \Delta_{(i)} \right)^{(m-k)} \xi(k; \Omega_{(i)}). \end{aligned}$$

• These lead to *exact* expressions of the risk functions of $\hat{\sigma}_{p,MLE}^2$, and of $\hat{\sigma}_{p,MRE}^2$ and $\hat{\sigma}_{p,ALB}^2$.

• When p = 2, expressions simplify. The mean becomes

$$EpMLE(\Delta) = 1 - \left(\frac{2}{\sqrt{n}}|\Delta|\right) \left\{ \phi\left(\frac{\sqrt{n}}{2}|\Delta|\right) - \left(\frac{\sqrt{n}}{2}|\Delta|\right) \left[1 - \Phi\left(\frac{\sqrt{n}}{2}|\Delta|\right)\right] \right\}.$$

- Follows from this that $\hat{\sigma}_{p,MLE}^2$ is negatively biased for σ^2 .
- **Question:** What happens when the number of sub-models increases indefinitely?
- **Theorem:** With n > 1 fixed, if as $p \to \infty$, $\max_{2 \le i \le p} |\Delta_{(i)} \Delta_{(i-1)}| \to 0$, $\Delta_{(1)} \to -\infty$, and $\Delta_{(p)} \to \infty$, then (i) Eff $\left(\hat{\sigma}_{p,MLE}^2 : \hat{\sigma}_{UMVU}^2\right) \to \frac{2n^2}{(n-1)(2n-1)} > 1$; (ii) Eff $\left(\hat{\sigma}_{p,MRE}^2 : \hat{\sigma}_{UMVU}^2\right) \to \frac{2(n+2)^2}{(n-1)(2n+7)} > 1$; (iii) Eff $\left(\hat{\sigma}_{p,MRE}^2 : \hat{\sigma}_{p,MLE}^2\right) \to \frac{(2n-1)(n+2)^2}{(2n+7)n^2} > 1$; and (iv) Eff $\left(\hat{\sigma}_{p,MRE}^2 : \hat{\sigma}_{MRE}^2\right) \to \frac{2(n+2)^2}{(n+1)(2n+7)} < 1$. Also, in the limit, $\hat{\sigma}_{p,ALB}^2$ is dominated by $\hat{\sigma}_{UMVU}^2$.
- From (iv), the advantage of exploiting \mathcal{M}_p could be *lost* forever when p increases!

Representation: Limiting Bayes Estimator

• **Theorem:** Under \mathcal{M}_p with μ_{i_0} the true mean,

$$\frac{\hat{\sigma}_{p,LB}^2}{\sigma^2} \stackrel{d}{=} \frac{W}{n} \left\{ 1 + H(\boldsymbol{T}) \right\},\,$$

where

$$\boldsymbol{T} = (T_1, T_2, \dots, T_p)' = \boldsymbol{V} / \sqrt{W},$$
$$H(\boldsymbol{T}) = \sum_{i=1}^p \theta_i(\boldsymbol{T}) T_i^2;$$
$$\theta_i(\boldsymbol{T}) = \frac{(1+T_i^2)^{-(n/2)}}{\sum_{j=1}^p (1+T_j^2)^{-(n/2)}}, \quad i = 1, 2, \dots, p$$

- However, even with this nice-looking representation, it is difficult to obtain exact expressions for the mean and variance.
- Developed 2nd-order approximations, but were not so satisfactory when $n \leq 15$.
- In the comparisons, we resorted to simulations to approximate the risk function of $\hat{\sigma}_{p,LB}^2$.

Table 1: 2nd-order approximation and simulation results for the mean and variance functions of $\hat{\sigma}_{p,LB}^2/\sigma^2$, and the risk function of $\hat{\sigma}_{p,LB}^2$ for different combinations of p, Δ , and n. For each combination, 10000 simulation replications were performed.

Combinations		Mean		Variance		Risk	
of p and $\boldsymbol{\Delta}$	n	Appr.	Sim.	Appr.	Sim.	Appr.	Sim.
$\Delta = (-0.25, 0, 0.25)$ p=3	$3 \\ 10 \\ 30$	$0.802 \\ 0.956 \\ 0.988$	$0.961 \\ 0.989 \\ 0.994$	$0.450 \\ 0.182 \\ 0.066$	$0.650 \\ 0.204 \\ 0.067$	$0.489 \\ 0.184 \\ 0.066$	$0.651 \\ 0.204 \\ 0.067$
$\Delta = (-0.5, 0, 0.5)$ = 3	$3 \\ 10 \\ 30$	$0.844 \\ 0.994 \\ 1.024$	$0.937 \\ 0.985 \\ 1.004$	$0.401 \\ 0.204 \\ 0.074$	$0.652 \\ 0.206 \\ 0.067$	$0.425 \\ 0.204 \\ 0.074$	$0.656 \\ 0.206 \\ 0.067$
$\Delta = (0, 0.25, 0.50)$ p=3	$\begin{array}{c} 3\\10\\30\end{array}$	$0.951 \\ 1.005 \\ 1.004$	$1.036 \\ 1.006 \\ 1.004$	$0.596 \\ 0.201 \\ 0.068$	$0.745 \\ 0.205 \\ 0.068$	$0.598 \\ 0.201 \\ 0.068$	$0.746 \\ 0.205 \\ 0.068$
$\Delta = (0, 0.5, 1)$ p=3	$3 \\ 10 \\ 30$	$1.071 \\ 1.025 \\ 1.012$	$1.106 \\ 1.021 \\ 1.004$	$0.726 \\ 0.220 \\ 0.070$	$0.886 \\ 0.220 \\ 0.068$	$0.731 \\ 0.220 \\ 0.071$	$0.898 \\ 0.221 \\ 0.068$
$\Delta = (-0.25:0.0625:0.25)$	$\begin{array}{c} 3\\10\\30\end{array}$	$0.857 \\ 0.970 \\ 0.986$	0.977 0.984 0.993	$0.495 \\ 0.185 \\ 0.065$	$0.648 \\ 0.203 \\ 0.066$	$0.515 \\ 0.186 \\ 0.065$	$0.649 \\ 0.204 \\ 0.066$
$ \Delta = (-0.25:0.03125:0.25) $ $ p = 17 $	$\begin{array}{c} 3\\10\\30\end{array}$	$0.867 \\ 0.972 \\ 0.987$	$0.974 \\ 0.988 \\ 0.994$	$0.504 \\ 0.186 \\ 0.065$	$0.670 \\ 0.202 \\ 0.067$	$0.522 \\ 0.187 \\ 0.065$	$0.671 \\ 0.202 \\ 0.067$
$\Delta = (0:0.0625:0.5)$ p=9	3 10 30	$\begin{array}{c} 0.992 \\ 1.021 \\ 1.013 \end{array}$	$1.031 \\ 1.033 \\ 1.013$	$0.642 \\ 0.206 \\ 0.069$	$\begin{array}{c} 0.742 \\ 0.221 \\ 0.071 \end{array}$	$0.642 \\ 0.206 \\ 0.069$	$\begin{array}{c} 0.743 \\ 0.222 \\ 0.071 \end{array}$
$\Delta = (0:0.03125:0.5)$ p=17	3 10 30	$1.000 \\ 1.024 \\ 1.015$	$1.042 \\ 1.030 \\ 1.018$	$0.652 \\ 0.207 \\ 0.069$	$0.755 \\ 0.215 \\ 0.071$	$0.652 \\ 0.207 \\ 0.069$	$0.756 \\ 0.216 \\ 0.071$



Figure 1: Relative efficiencies of σ^2 estimators for p = 2 and $\Delta = (0, \Delta)$ for $n \in \{3, 10, 30\}$. For the (p,LB) the (connected) scatterplot represents the simulated estimates of relative efficiency based on the 10000 replications for each Δ .

Combinations	Efficiency %						
of p and $\boldsymbol{\Delta}$	n	UMVU	MVU MRE LB		pMLE	pMRE	ALB
$\Delta = (-0.25, 0, 0.25)$ p=3	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 156 \ (148, \ 165) \\ 110 \ (110, \ 117) \\ 105 \ (102, \ 107) \end{array}$	$173 \\ 117 \\ 105$	$222 \\ 123 \\ 107$	$\begin{array}{c} 14\\71\\90\end{array}$
$\Delta = (-0.5, 0, 0.5)$ p=3	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 162 \ (153, \ 168) \\ 107 \ (107, \ 114) \\ 97 \ (99, \ 105) \end{array}$	$183 \\ 119 \\ 106$	$209 \\ 124 \\ 110$	17 72 88
$\Delta = (0, 0.25, 0.50)$ p=3	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 137 \ (133, 146) \\ 103 \ (102, 109) \\ 98 \ (99, 105) \end{array}$	$164 \\ 114 \\ 104$	$226 \\ 126 \\ 108$	$12 \\ 65 \\ 87$
$\Delta = (0, 0.5, 1)$ p=3	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 109 \ (111, \ 125) \\ 98 \ (97, \ 105) \\ 102 \ (100, \ 106) \end{array}$	$165 \\ 114 \\ 104$	221 128 110	$13 \\ 65 \\ 86$
	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 154 \ (148, \ 163) \\ 109 \ (110, \ 118) \\ 102 \ (102, \ 108) \end{array}$	$173 \\ 117 \\ 105$	$221 \\ 122 \\ 105$	$14 \\ 71 \\ 91$
$\Delta = (-0.25:2^{-5}:0.25)$ p=17	$3 \\ 10 \\ 30$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$150 (148, 161) \\ 113 (109, 118) \\ 104 (102, 108)$	$173 \\ 116 \\ 105$	$221 \\ 122 \\ 105$	$14 \\ 71 \\ 91$
$\Delta = (0:2^{-4}:0.5)$ p=9	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 136 \ (132, 146) \\ 100 \ (100, 108) \\ 98 \ (97, 103) \end{array}$	$163 \\ 114 \\ 104$	$225 \\ 125 \\ 107$	$\begin{array}{c} 12 \\ 65 \\ 87 \end{array}$
$\Delta = (0:2^{-5}:0.5)$ p=17	${3 \\ 10 \\ 30}$	$100 \\ 100 \\ 100$	$200 \\ 122 \\ 106$	$\begin{array}{c} 134 \ (132, 146) \\ 102 \ (100, 107) \\ 99 \ (96, 103) \end{array}$	$163 \\ 114 \\ 104$	$225 \\ 125 \\ 107$	$12 \\ 65 \\ 87$

Table 2: Relative efficiencies of the variance estimators for different combinations of p, Δ , and n. For the (p,LB) estimator, 95% empirical confidence intervals are shown in parenthesis.



A contour plot as a function of p and deltamax, symmetric case

A contour plot as a function of p and deltamax, asymmetric case



Figure 2: Relative efficiencies of $\hat{\sigma}_{p,MRE}^2$ (pMRE) wrt $\hat{\sigma}_{MRE}^2$ (MRE) in a symmetric and asymmetric Δ cases, as a function of Δ_{\max} and p for sample size of n = 10. Symmetric case of form $\mathbf{\Delta} = [-\Delta_{max} : \Delta_{max}/(p-1) : \Delta_{max}]$; while asymmetric case of form $\mathbf{\Delta} = [0 : \Delta_{max}/(2(p-1)) : \Delta_{max}]$.

7. Recap: $\tau(t) = F(t)$ Estimators

• \mathcal{M} -UMVU:

$$z_1(t) = \frac{\sqrt{n}}{n-1} \left(\frac{t-\bar{X}}{S} \right)$$

$$\hat{\tau}_{UMVU}(t) = \mathcal{T}\left(\frac{\sqrt{n-2}z_1(t)}{\sqrt{1-z_1(t)^2}}; n-2\right) I\{|z_1(t)| \le 1\} + I\{z_1(t) > 1\}$$

• \mathcal{M}_p -UMVU:

$$z_{3i}(t) = \frac{1}{\sqrt{n}} \left(\frac{t - \mu_i}{\hat{\sigma}_i} \right), \quad i = 1, 2, \dots, p$$

$$\hat{\tau}_{UMVU,i}(t) = \mathcal{T}\left(\frac{\sqrt{n-1}z_{3i}(t)}{\sqrt{1-z_{3i}(t)^2}}; n-1\right) I\{|z_{3i}(t)| \le 1\} + I\{z_{3i}(t) > 1\}$$

$$\hat{\tau}_{p,UMVU}(t) = \hat{\tau}_{UMVU,\hat{M}}(t) = \sum_{i=1}^{p} I\{\hat{M} = i\}\hat{\tau}_{UMVU,i}(t)$$

• \mathcal{M}_p -Limiting Bayes:

$$\hat{\tau}_{p,LB}(t) = \sum_{i=1}^{p} \left\{ \frac{\left(\hat{\sigma}_{i}^{2}\right)^{-n/2}}{\sum_{j=1}^{p} \left(\hat{\sigma}_{j}^{2}\right)^{-n/2}} \right\}$$

• Also obtained distributional representations for these estimators which show that their distributions depend on (t, μ, σ^2) only through $(\xi(t), \Delta)$, where

$$\xi(t) = \frac{t - \mu_{i_0}}{\sigma} =$$
standardized *t*-value.

- But, even with the representations, *no* exact expressions of their risk functions were obtained.
- Comparisons for the $\tau(t)$ -estimators were therefore performed through simulations.

Results of Comparisons

- For the τ -estimators, efficiencies are relative to $\hat{\tau}_{UMVU}$.
- To compare globally the τ -estimators, we approximated the risk functions arising from the global loss function

$$L_3(\hat{\tau},(\mu,\sigma^2)) = \int \left[\hat{\tau}(t) - \Phi\left(\frac{t-\mu}{\sigma}\right)\right]^2 \Phi\left(\frac{dt-\mu}{\sigma}\right).$$



Figure 3: Pointwise biases, variances, risks, and relative efficiencies of the four distribution estimators $\hat{\tau}$, for $\Delta = (-1, 0, 1)$ and sample size n = 10. For each standardized time point $\xi(t)$, 10000 simulation replications were performed.

Table 3: Relative global efficiencies (rel. to the UMVU estimator $\hat{\tau}_{UMVU}^2$) of the three distribution estimators for different combinations of p, Δ , and n. 10000 simulation replications were performed for each combination.

Combinations of p and $\boldsymbol{\Delta}$	n	$\substack{\text{pUMVU}\\\text{Eff }\%}$	LB Eff %	$\begin{array}{c} \mathrm{ALB} \\ \mathrm{Eff} \ \% \end{array}$	Combinations of p and Δ	$_{\rm Eff~\%}^{\rm pUMVU}$	LB Eff %	ALB Eff %
$\Delta = (-0.25, 0, 0.25)$ p=3	$3 \\ 10 \\ 30$	$316 \\ 190 \\ 107$	$704 \\ 419 \\ 194$	$395 \\ 197 \\ 107$	$\Delta = (-1, 0, 1)$ p=3	$105 \\ 108 \\ 352$	181 138 389	$117 \\ 111 \\ 359$
$\Delta = (-0.5, 0, 0.5)$ p=3	$3 \\ 10 \\ 30$	$169 \\ 94 \\ 83$	$367 \\ 158 \\ 110$	$194 \\ 96 \\ 84$	$\Delta = (-1:2^{-1}:1)$ p=5	109 87 82	199 115 108	123 89 82
$\Delta = (0, 0.25, 0.50)$ p=3	$3 \\ 10 \\ 30$	$274 \\ 190 \\ 158$	$428 \\ 220 \\ 163$	$334 \\ 197 \\ 159$	$\Delta = (-1:2^{-2}:1)$ p=9	$110 \\ 95 \\ 88$	$214 \\ 120 \\ 103$	125 98 89
$\Delta = (0, 0.5, 1)$ p=3	$3 \\ 10 \\ 30$	$180 \\ 152 \\ 145$	$239 \\ 166 \\ 181$	$209 \\ 157 \\ 146$	$\Delta = (-1:2^{-3}:1)$ p=17	$ \begin{array}{r} 111 \\ 98 \\ 97 \end{array} $	$227 \\ 123 \\ 103$	125 101 97
$\begin{array}{c} \mathbf{\Delta} = (-0.25:2^{-4}:0.25) \\ p = 9 \end{array}$	$3 \\ 10 \\ 30$	$321 \\ 205 \\ 128$	$779 \\ 530 \\ 288$	$403 \\ 213 \\ 129$	$\Delta = (-1:2^{-4}:1)$ p=33	$110 \\ 99 \\ 99 \\ 99$	$230 \\ 125 \\ 103$	$124 \\ 102 \\ 99$
$ \Delta = (-0.25:2^{-5}:0.25) $ p=17	$3 \\ 10 \\ 30$	$319 \\ 204 \\ 130$	$774 \\ 560 \\ 313$	$400 \\ 212 \\ 131$	$\Delta = (-1:2^{-5}:1)$ p=65	$111 \\ 99 \\ 99 \\ 99$	239 126 103	$125 \\ 102 \\ 100$
$\Delta = (0:2^{-4}:0.5)$ p=9	$3 \\ 10 \\ 30$	$274 \\ 196 \\ 174$	$431 \\ 204 \\ 131$	$334 \\ 204 \\ 175$	$\Delta = (-1:2^{-6}:1)$ p=129	$110 \\ 99 \\ 99 \\ 99$	$237 \\ 127 \\ 103$	$124 \\ 102 \\ 100$
$\Delta = (0:2^{-5}:0.5)$ p=17	$ \begin{array}{c} 3 \\ 10 \\ 30 \end{array} $	$276 \\ 203 \\ 176$	$432 \\ 209 \\ 125$	$336 \\ 211 \\ 178$	$\Delta = (-1:2^{-7}:1)$ p=257	111 99 99	$238 \\ 127 \\ 103$	$125 \\ 102 \\ 100$
$\begin{array}{c} \boldsymbol{\Delta} = (0, 1) \\ p = 2 \end{array}$	$ \begin{array}{c} 3 \\ 10 \\ 30 \end{array} $	$176 \\ 186 \\ 525$	$259 \\ 227 \\ 553$	$204 \\ 193 \\ 539$	$\Delta = (-1:2^{-8}:1)$ p=513	111 99 99	239 127 103	$125 \\ 102 \\ 100$



A contour plot as a function of n and deltamax for p=2

A contour plot as a function of n and deltamax for p=3



Figure 4: Relative global efficiencies of pLB with respect to UMVU in a asymmetric (p = 2) and symmetric (p = 3) Δ cases, as a function of Δ_{\max} and sample size n. Scenario 1: Corresponds to an asymmetric case of form $\Delta = (0, \Delta_{\max})$. Scenario 2: Corresponds to a symmetric case of form $\Delta = (-\Delta_{\max}, 0, \Delta_{\max})$. For each combination of (n, Δ) , 10000 simulation replications were performed.

8. Concluding Remarks

- In models with sub-models, and interest is to infer about a common parameter, possible approaches are:
- Approach I: Utilize procedures for a wider model subsuming the sub-models. Could lead to loss of efficiency.
- Approach II: Utilize a two-step approach: First step is to select the sub-model using the data; second step is to use a procedure (e.g., estimator) for the chosen sub-model, again using the same data.

Should recognize that properties of the two-step procedure will be different from the sub-model properties of the procedures.

• Approach III: Utilize a Bayesian framework. Assign a prior to the sub-models, and (conditional) priors on the parameters within the sub-models.

Resulting procedure is an adaptively weighted combination of the (Bayes) procedures in the sub-models.

- Approaches (II) and (III) appear preferable over approach (I), but when the number of sub-models is large, approach (I) may provide better estimators and a simpler determination of the properties.
- Hard to conclude which of approaches (II) or (III) is preferable. In the Gaussian model considered, approach (II) performed better in estimating the variance σ², but approach (III) performed better in estimating the distribution function.

This calls for further studies in more complicated settings, such as those that motivated our study.

• To conclude,

Observe Caution!

when doing inference after model selection especially when *double-dipping* on the data!