

Version of 02–04–10

## ANOMALIES IN RIDGE REGRESSION: ADDENDUM

DONALD R. JENSEN AND DONALD E. RAMIREZ

### 1. INTRODUCTION

We wish to thank Professors Kapat and Goel (2010) (hereafter KG(2010)) for clarifying the noninvariance of ridge regression to the choice of parametrization, a concern voiced by Smith and Campbell (1980), their discussants, and others. KG(2010) focus mainly on one aspect of our “Foundations” paper (hereafter JR(2008)), namely, constrained optimization in linear inference. We mostly concede their corrections, noting that essentials are found in Marquardt (1963), Meeter (1966), Hoerl and Kennard (1970), and recently Davidov (2006). However, the principal anomaly of JR(2008) remains: The theory and practice of ridge regression do *not* rest on constrained optimization.

Here we trace conventions for taking  $\{\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}\}$ , in original coordinates, through centering and scaling into  $\{\mathbf{Y} = \mathbf{Z}\boldsymbol{\theta} + \boldsymbol{\epsilon}\}$  with  $\mathbf{A} = \mathbf{Z}'\mathbf{Z}$  in “correlation form”, the model  $M_R$  of KG(2010), but now distinguishing  $\boldsymbol{\theta}$  from  $\boldsymbol{\beta}$ . In practice, the original scale typically is of overriding concern to the analyst; hence the routine mapping back to that scale. In fact, an algorithm in Section 4.2 of JR(2008) shows connections between  $\{\mathbf{Y} = \mathbf{Z}\boldsymbol{\theta} + \boldsymbol{\epsilon}\}$ , where solutions diminish monotonically in length, and lengths of projections back onto the original scale, which need not be monotone.

### 2. STANDARDIZATION

Constrained optimization in regression traces to Marquardt (1963), seeking a compromise between Gauss–Newton and gradient methods in nonlinear estimation, and solving iteratively through locally linear models. Noting that gradient methods are not scale–invariant, Marquardt (1963; p.436) for convenience scales the  $\boldsymbol{\beta}$ –space using standard deviations of the columns of  $\mathbf{X}$ , transforming elements of  $\mathbf{X}'\mathbf{X}$  into “simple correlation coefficients,” as used widely in least–squares problems “for improving the numerical aspects of computing procedures.” As in Belsley (1986), the latter includes the conditioning of linear systems and stability of their solutions. This standardization pervades much of the ridge regression literature, beginning with Hoerl and Kennard (1970) and onward.

That solutions are mapped routinely back onto the original scale is seen in the *Ridge* option of *Proc Reg* in the *SAS* system, as used in Myers (1990) for the Hospital Manpower Data and as verified in JR(2008). An antecedent is Marquardt’s (1963) algorithm: At each iteration the local linear model is transformed into “correlation form”; the constrained solution is found; and this is projected back onto the original scale before proceeding to the next iterate. Letting  $\mathbf{D}_s = \text{Diag}(s_1^{-1}, \dots, s_p^{-1})$  as in KG(2010), the constraint  $\{\boldsymbol{\theta}'\boldsymbol{\theta} =$

$c^2$  maps back to  $\{\beta' D_s^{-2} \beta = c^2\}$ , often counter-intuitive to the user and subject to the vagaries of typically irrelevant lengths of the columns of  $\mathbf{X}$ . That ridge solutions  $\|\widehat{\boldsymbol{\theta}}_{R_k}\|^2 = \widehat{\boldsymbol{\beta}}_{R_k}' D_s^{-2} \widehat{\boldsymbol{\beta}}_{R_k}$  decrease monotonically in  $k$  is an artifact of an algorithm chosen for numerical stability. To the contrary, users instead may insist that  $\|\widehat{\boldsymbol{\beta}}_{R_k}\|$  is essential, and then seek the value  $\widehat{\boldsymbol{\beta}}_R(\widehat{k})$  so as (i) to achieve a natural constraint  $\{\|\widehat{\boldsymbol{\beta}}_{R_k}\| = c\}$ , and (ii) to minimize the residual sum of squares on that scale. To these ends, the algorithm in Section 4.2 of JR(2008) bridges this gap, demonstrating explicitly the manner in which solutions from  $\{\mathbf{Y} = \mathbf{Z}\boldsymbol{\theta} + \boldsymbol{\epsilon}\}$ , when projected as lengths onto the original scale, will be minimizing there. Specifically, for fixed  $c$  again define the set

$$\Lambda(c) = \{k : \|\widehat{\boldsymbol{\beta}}_R(k)\| = c\},$$

and let  $k_c = \min\{\Lambda(c)\}$ . Then  $\widehat{\boldsymbol{\beta}}_R(k_c)$  achieves its designated constraint and is minimizing, all on the original scale. KG(2010) claim that  $\Lambda(c)$  is either a singleton set or is empty. While true for solutions in “correlation form,” the cardinality of  $\Lambda(c)$  may increase when projected back onto the original scale, as noted and illustrated correctly in JR(2008).

### 3. CONCLUSIONS

Pretensions to the contrary, ridge regression remains ill-posed mathematically. Under objective constraints, either  $\{\beta'\beta = c^2\}$  or  $\{\beta'\beta \leq c^2\}$  in a generic model  $\{\mathbf{Y} = \mathbf{X}\beta + \boldsymbol{\epsilon}\}$ , constrained optimization returns a solution  $[\widehat{\boldsymbol{\beta}}_R(\widehat{k}), \widehat{k}, c^2]$ . That an objective  $c^2$  typically is missing, is acknowledged in the myriad choices for  $k$  advocated in the literature, including the early “ridge trace” of Hoerl and Kennard (1970).

The overriding determinant of the need for ridge regression, or other biased estimation, is conditioning of the system  $\{\mathbf{X}'\mathbf{X}\beta = \mathbf{X}'\mathbf{Y}\}$ . Possibilities to improve the conditioning of  $\mathbf{X}'\mathbf{X}$  include  $\{\mathbf{X}'\mathbf{X} \rightarrow (\mathbf{X}'\mathbf{X} + \mathbf{B})\}$  with  $\mathbf{B}$  positive definite. For  $k > 0$  it is known (Riley (1955)) that  $(\mathbf{X}'\mathbf{X} + k\mathbf{I}_p)$  is better conditioned than  $\mathbf{X}'\mathbf{X}$ , cited as justification for ridge regression in Marshall and Olkin (1970). Nonetheless, constrained optimization is not a panacea: Given any “ridge-type” system  $\{(\mathbf{X}'\mathbf{X} + k\boldsymbol{\Delta})\beta = \mathbf{X}'\mathbf{Y}\}$ , under generalized constraints as in Section 1 of KG(2010) and McDonald (1982), it remains open (apart from  $\boldsymbol{\Delta} = k\mathbf{I}_p$ ) as to circumstances where its conditioning might be enhanced over the original. If not, this could further undermine conditioning of the system, and exacerbate the problem through an inappropriate and counter-productive choice of constrained optimization.

The focus of KG(2010) on constrained optimization, and our continuation in this reply, are largely beside the point, distracting, and obscuring from the main thesis of JR(2008): Whatever its niceties, constrained optimization has essentially nothing to do with ridge regression as practiced. Numerous procedures have been advanced for determining  $k$ , to include ridge traces, choices assuring dominance in mean square over Ordinary Least Squares (*OLS*) or, *inter se*, among other biased estimators; and other approaches. In particular, variances, biases, and mean square errors, as given in Section 4 of Hoerl and Kennard (1970); and numerous other properties of ridge solutions adopted by generations hence: These all assert incorrectly that  $\widehat{\boldsymbol{\beta}}_{R_k}$ , intended as a constrained estimator, is linear in the

*OLS* solution  $\hat{\beta}_L$ . Granted, the ridge system  $\{(\mathbf{X}'\mathbf{X} + k\mathbf{I}_p)\hat{\beta}_{R_k} = \mathbf{X}'\mathbf{Y}\}$  does derive from initial steps of constrained optimization. However, this in itself fails to justify repeated claims that ridge regression rests on constrained optimization. To be correct, the final step must deliver the constrained solution. As noted earlier, conditioning arguments alone deliver the identical system  $\{(\mathbf{X}'\mathbf{X} + k\mathbf{I}_p)\hat{\beta}_{R_k} = \mathbf{X}'\mathbf{Y}\}$  without recourse to constraints. In short, once again, estimators constrained to the  $c$ -sphere, or to the  $c$ -ball in Euclidean  $p$ -space, cannot be linear in  $\hat{\beta}_L$ , as we emphasized from several perspectives in JR(2008). Specifically, a solution constrained to a ball of radius  $c$  *could* be unbiased, depending on how its probability is concentrated there. Nothing offered by KG(2010) can alter these facts, totally debilitating as they are to claims that constrained optimization is the foundation for ridge regression as known and practiced.

Ridge regression having failed its vaunted credentials, we introduced in JR(2008) *surrogate estimators* as modifications of ridge based on conditioning, and we showed numerically in our case study that, in contrast to ridge, surrogate estimators possess desirable monotone properties. As the ridge parameter  $k$  evolves, we since have shown in theory that ridge estimators typically exhibit erratic divergence from those of orthogonal systems, often reverting back to *OLS* in the limit. In contrast, surrogate solutions are seen to converge monotonically to those from orthogonal systems. This work, titled “Surrogate Models in Ill-Conditioned Systems”, is scheduled to appear in a forthcoming issue of the *Journal of Statistical Planning and Inference*.

#### REFERENCES

- [1] **Note:** Items cited but not listed are found in JR(2008) and KG(2010).
- [2] Belsley, D.A. (1986). “Centering, the constant, first-differencing, and assessing conditioning.” Chapter 5 in *Model Reliability*, D.A. Belsley and E. Kuh, eds., pp. 117–152. MIT Press, Cambridge, MA.

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF VIRGINIA, CHARLOTTESVILLE, VA 22904–4137  
*E-mail address:* djensen@vt.edu; der@virginia.edu