

## DISCUSSION OF: A STATISTICAL ANALYSIS OF MULTIPLE TEMPERATURE PROXIES: ARE RECONSTRUCTIONS OF SURFACE TEMPERATURES OVER THE LAST 1000 YEARS RELIABLE?<sup>1</sup>

BY ALEXEY KAPLAN<sup>2</sup>

*Lamont–Doherty Earth Observatory of Columbia University*

McShane and Wyner (2011) (hereinafter MW2011) demonstrated that in many cases a comprehensive data set of  $p = 1138$  proxies [Mann et al. (2008)] did not predict Northern Hemisphere (NH) mean temperatures significantly better than random numbers. This fact is not very surprising in itself: the unsupervised selection of good predictors from a set of  $p \gg n$  proxies of varying sensitivities might be too challenging a task for any statistical method ( $p/n_c \approx 10$ ; only  $n_c = 119$  out of total  $n = 149$  years were used for calibration in MW2011 cross-validated reconstructions). However, some types of noise<sup>3</sup> systematically outperformed the real proxies (see two bottom panels of MW2011, Figure 10). This finding begs further investigation: what do these random numbers have that real proxies do not?

To investigate this question, the present analysis uses ridge regression [RR, Hoerl and Kennard (1970)] instead of the Lasso [Tibshirani (1996)].<sup>4</sup> The regression model used by MW2011 with Lasso and here with RR is

$$y = X\beta + \beta_0\mathbb{1}_n + \varepsilon,$$

where  $y$  is a column vector of  $n$  observations (annual NH temperatures),  $\varepsilon$  is random error,  $X$  is a known  $n \times p$  matrix of predictors (climate proxies). A vector of regression coefficients  $\beta$  and an intercept constant  $\beta_0$  are to be determined. A column  $n$ -vector  $\mathbb{1}_n$  has all components equal one. Proxy records are standardized before use; in cross-validation experiments standardization is repeated for each calibration period.

---

Received October 2010; revised November 2010.

<sup>1</sup>Lamont–Doherty Earth Observatory contribution number 7438.

<sup>2</sup>Supported by grants from the NSF (ATM-0902436), NOAA (NA07OAR4310060), and NASA (NNX09AF44G).

*Key words and phrases.* Paleoclimate, statistical climate reconstructions, cross-validation, ridge regression, autoregressive processes, kriging.

<sup>3</sup>Pseudoproxies used by MW2011 are called “noise” here; in climate research, pseudoproxies are synthetic combinations of a climate signal with some noise; without the former, it is a pure noise.

<sup>4</sup>The difference is in the penalty norm: Lasso uses  $L_1$  while RR uses  $L_2$ . MW2011 have also argued that a rough performance similarity should exist between different methods for  $p \gg n$  problems.

Let  $w$  be a column  $n_c$ -vector such that  $w^T \mathbb{1}_{n_c} = 1$ . Define matrix-valued functions  $\mathcal{W}[w] = I - \mathbb{1}_{n_c} w^T$  and  $\mathcal{R}[S, \lambda, w] = S_{vc}(S_{cc} + \lambda I)^{-1} \mathcal{W}[w] + \mathbb{1}_{n_v} w^T$ , where  $S$  is a positive semidefinite  $n \times n$  matrix,  $\lambda > 0$  is the ridge parameter found as a minimizer of the generalized cross-validation function [GCV, Golub et al. (1979)], matrix (or vector) subscripts  $c$  or  $v$  hereinafter indicate submatrices corresponding to the calibration or validation periods, respectively. The RR reconstruction  $\hat{y}_v$  of temperatures in the validation period (a “holdout block” of  $n_v = 30$  consecutive years) is a linear transformation:  $\hat{y}_v = \mathcal{R}[S_p, \lambda, e] y_c$ , where  $S_p = \tilde{X} \tilde{X}^T / p$ ,  $\tilde{X}$  is the standardized version of  $X$ , and  $e = n_c^{-1} \mathbb{1}_{n_c}$ .

Using these formulas, the RR version of the MW2011 cross-validation tests were performed for real proxies and for some noise types. Results are shown in Figure 1. The cross-validated root mean square error (RMSE) of the RR reconstructions are smaller than Lasso values (cf. MW2011, Figure 9), but the relative performance in different experiments appears consistent between RR and Lasso. As in the Lasso case, noise with high temporal persistence, that is, simulated by the Brownian motion or by the first-order autoregressive process AR(1) with a parameter  $\varphi \geq 0.9$ , outperformed proxies. Figure 2 illustrates the time dependence of

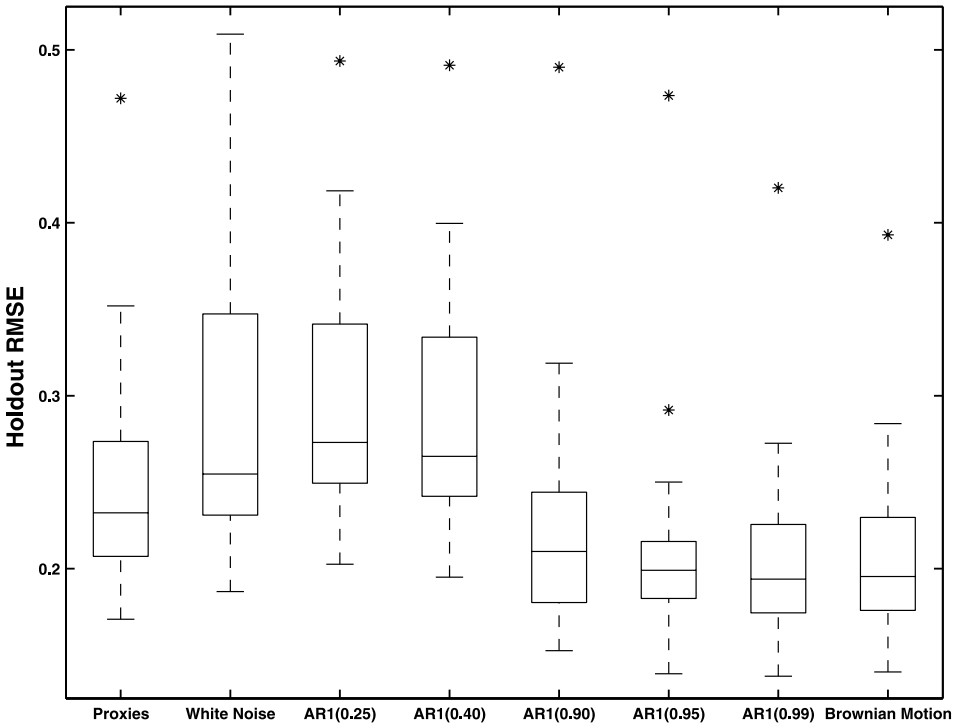


FIG. 1. Cross-validated RMSE on 120 30-year holdout blocks for the RR reconstructions from real climate proxies and from the random noise (one realization for each noise experiment); cf. MW2011, Figure 9.

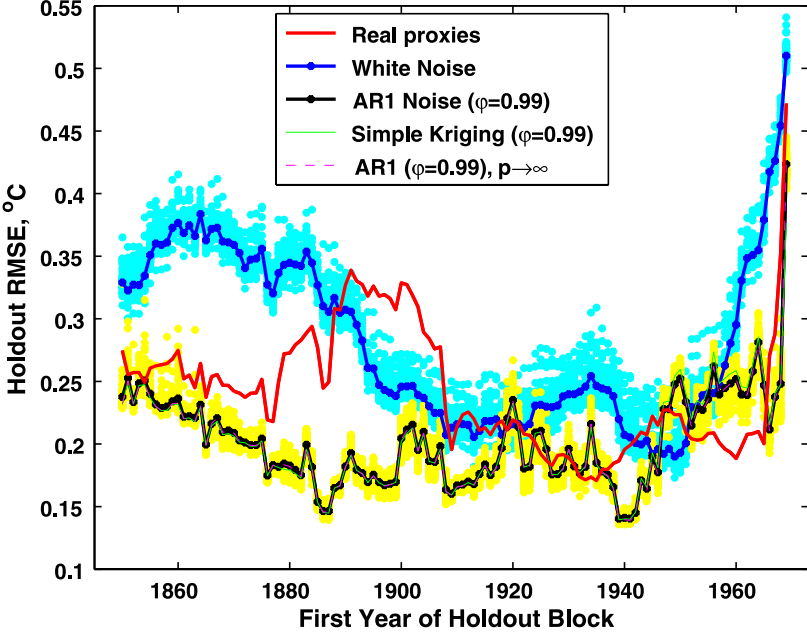


FIG. 2. Holdout RMSE for RR reconstructions as a function of time for real proxies (red) and two 100-member ensemble means: white noise (blue) and AR(1) noise with  $\varphi = 0.99$  (black). The probability limit ( $p \rightarrow \infty$ ) for the latter is shown by magenta dashes. Holdout RMSE for simple kriging of the NH mean temperature index using an exponential semivariogram [Le and Zidek (2006)]  $\gamma(\tau) = \lambda_{\min} + 1 - \exp[\tau \ln \varphi]$  with the GCV-selected nugget  $\lambda_{\min} = \ell(\Phi, 0)$  and long decorrelation scale  $-1/\ln(\varphi) = 99.5$  years ( $\tau$  is time in years) is shown by the green line. Individual ensemble members are shown by magenta and yellow dots, respectively.

the holdout error for the real-proxy, white-noise, and  $\varphi = 0.99$  AR(1) cases. There is a general similarity between these and the corresponding curves in Figure 10 by MW2011.

Note that a traditional approach to hypothesis testing would evaluate an RMSE corresponding to a regression of temperature data ( $y$ ) on real proxies ( $X$ ) in the context of the RMSE probability distribution induced by the assumed distribution of  $y$  under the hypothesized condition (e.g.,  $\beta = 0$ ). However, MW2011 evaluate the RMSE of real proxies in the context of the RMSE distribution induced by random values in  $X$ , not  $y$ . Such an approach to testing a null hypothesis would be appropriate for an inverse relationship, that is,  $X = y\beta^T + \mathbb{1}_n\beta_0^T + \varepsilon$ . When used with a direct regression model here, however, it results in the RMSE distribution with a surprising feature: when  $p \rightarrow \infty$ , RMSE values for individual realizations of the noise matrix  $X$  converge in probability to a constant.

This convergence occurs because the columns  $x$  of  $X$  in the noise experiments are i.i.d. from the noise distribution; AR(1) with  $\varphi = 0.99$  is considered here:  $x \sim \mathcal{N}(0, \Phi)$ ,  $\Phi = (\varphi^{|i-j|})$ . The columns of  $\tilde{X}$  are i.i.d. too, hence the random matrix

$S_p = \tilde{X}\tilde{X}^T/p$  is an average of  $p$  i.i.d. variates  $\tilde{x}\tilde{x}^T$ . Expectation  $\Psi = \mathbf{E}\tilde{x}\tilde{x}^T$  exists; its elements are computed as expectations of ratios and first inverse moments of quadratic forms in normal variables [Jones (1986, 1987)]. The weak law of large numbers applies, so  $S_p \xrightarrow{P} \Psi$ . Since the GCV function depends on  $S$  and  $w$  as well as on  $\lambda$ , its minimizing  $\lambda$  will depend on these parameters too:  $\lambda_{\min} = \ell[S, w]$ . Here GCV is assumed well-behaved, so that  $\ell$  is a single-valued function, continuous at  $(\Psi, e)$ . From the definition of  $\mathcal{R}$ ,  $\mathcal{B}[S, e] \equiv \mathcal{R}[S, \ell[S, e], e]$  will also be continuous at  $S = \Psi$ , thus  $S_p \xrightarrow{P} \Psi$  implies  $\hat{y}_v = \mathcal{B}[S_p, e]_{y_c} \xrightarrow{P} \mathcal{B}[\Psi, e]_{y_c}$ .

When  $p$  is finite but large, like  $p = 1138$ , reconstructions based on individual realizations of a noise matrix  $X$  are dominated by their constant components, especially when  $\varphi \approx 1$ : note the small scatter of RMSE values in the ensemble of AR(1) with  $\varphi = 0.99$  (yellow dots in Figure 2). The probability limit  $\hat{y}_v = \mathcal{B}[\Psi, e]_{y_c}$  yields RMSE values (magenta dash in Figure 2) that are very close ( $1.3 \cdot 10^{-3} \text{ }^\circ\text{C}$  RMS difference) to the ensemble mean RMSE (black curve in Figure 2). To interpret this non-random reconstruction, consider its simpler analogue, using neither proxy standardization nor a regression intercept ( $\beta_0$ ). Then, if the assumptions on the GCV function change accordingly,  $\hat{y}_v \xrightarrow{P} \mathcal{B}[\Phi, 0]_{y_c} = \Phi_{vc}[\Phi_{cc} + \ell(\Phi, 0)I]^{-1}y_c$ , that is, a prediction of  $y_v$  from  $y_c$  by ‘‘simple kriging’’ [Stein (1999, page 8)], which in atmospheric sciences is called objective analysis or optimal interpolation [Gandin (1963)]. The RMSE corresponding to this solution is shown in Figure 2 (green line): it is quite close to the ensemble mean RMSE for AR(1) noise with  $\varphi = 0.99$  (RMS difference is  $5.4 \cdot 10^{-3} \text{ }^\circ\text{C}$ ). The solution  $\mathcal{B}[\Psi, e]_{y_c}$ , to which the noise reconstructions without simplifications converge as  $p \rightarrow \infty$ , is more difficult to interpret. Still, it has a structure of an objective analysis solution and gives results that are similar to simple kriging: the RMS difference between the two reconstructions over all holdout blocks is  $7.7 \cdot 10^{-3} \text{ }^\circ\text{C}$ .

Due to the large value of  $p$  in the MW2011 experiments, their tests with the noise in place of proxies essentially reconstruct holdout temperatures by a kriging-like procedure in the temporal dimension. The covariance for this reconstruction procedure is set by the temporal autocovariance of the noise. Long decorrelation scales ( $\varphi \geq 0.95$ ) gave very good results, implying that long-range correlation structures carry useful information about predictand time series that is not supplied by proxies. By using such a noise for their null hypothesis, MW2011 make one skillful model (multivariate linear regression on proxies) compete against another (statistical interpolation in time) and conclude that a loser is useless. Such an inference does not seem justified.

Modern analysis systems do not throw away observations simply because they are less skillful than other information sources: instead, they combine information. MW2011 experiments have shown that their multivariate regressions on the proxy data would benefit from additional constraints on the temporal variability of the target time series, for example, with an AR model. After proxies are combined with such a model, a test for a significance of their contributions to the common product could be performed.

**Acknowledgements.** Generous technical help and many useful comments from Jason Smerdon and very helpful presentation style guidance from Editor Michael Stein are gratefully acknowledged.

### SUPPLEMENTARY MATERIAL

**Data and codes** (DOI: [10.1214/10-AOAS398MSUPP](https://doi.org/10.1214/10-AOAS398MSUPP); .zip). This supplement contains a tar archive with all data files and codes (Matlab scripts) needed for reproducing results presented in this discussion. Dependencies between files in the archive and the order in which Matlab scripts have to be executed are described in the file *README\_final*, also included into the archive.

### REFERENCES

- GANDIN, L. S. (1963). *Objective Analysis of Meteorological Fields*. Gidrometeorologicheskoye Izdatel'stvo, Leningrad. Translated from Russian, Israeli Program for Scientific Translations. Jerusalem, 1965.
- GOLUB, G. H., HEATH, M. and WAHBA, G. (1979). Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics* **21** 215–223. [MR0533250](#)
- HOERL, A. E. and KENNARD, R. W. (1970). Ridge regression: Biased estimation for non-orthogonal problems. *Technometrics* **12** 55–67.
- JONES, M. C. (1986). Expressions for inverse moments of positive quadratic forms in normal variables. *Austral. J. Statist.* **28** 242–250. [MR0860469](#)
- JONES, M. C. (1987). On moments of ratios of quadratic forms in normal variables. *Statist. Probab. Lett.* **6** 129–136. [MR0907273](#)
- LE, N. D. and ZIDEK, J. V. (2006). *Statistical Analysis of Environmental Space–Time Processes*. Springer, New York. [MR2223933](#)
- MANN, M. E., ZHANG, Z., HUGHES, M. K., BRADLEY, R. S., MILLER, S. K., RUTHERFORD, S. and NI, F. (2008). Proxy-based reconstructions of hemispheric and global surface temperature variations over the past two millennia. *Proc. Natl. Acad. Sci. USA* **105** 13252–13257.
- MCSHANE, B. B. and WYNER, A. J. (2011). A statistical analysis of multiple temperature proxies: Are reconstructions of surface temperatures over the last 1000 years reliable? *Ann. Appl. Statist.* **5** 5–44.
- STEIN, M. L. (1999). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer, New York. [MR1697409](#)
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the Lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. [MR1379242](#)

LAMONT–DOHERTY EARTH OBSERVATORY  
 61 ROUTE 9W  
 P.O. BOX 1000  
 PALISADES, NEW YORK 10964  
 USA  
 E-MAIL: [alexeyk@ldeo.columbia.edu](mailto:alexeyk@ldeo.columbia.edu)  
 URL: <http://rainbow.ldeo.columbia.edu/~alexeyk>