

# Spatial correlation

Hans Martinez

2023-07-19

## 1 Replication (updated)

The results of the third try in replicating Table 1 of Müller and Watson (2022b) is displayed in Table 1. In all models, locations are  $s_l \sim U(0, 1)$  and  $\beta_0 = 0$ . In all models  $e_l \sim \mathcal{G}_{exp}(c_{0.030})$ , except model 2. In model 1,  $x_l = 1$ . In model 2,  $x_l \sim \mathcal{G}_{exp}(c_{0.03}/2)$  and  $e_l \sim \mathcal{G}_{exp}(c_{0.03}/2)$ . In model 3,  $x_l$  is a step function with  $x_l = -0.15$  for the 85% of the locations closest to  $s = 0$ , and  $x_l = 0.85$  for the remaining locations closest to  $s = 1$ . In model 4,  $x_l$  follows a demeaned random walk.

I ran 1000 simulations for each model. I used 250 observations for each iteration. One location per observation. Locations were fixed. The value of  $\beta$  was 0 and  $\bar{\rho}$ , 0.0301.

The authors do not specify clustering at all. The rejection frequencies are the percentage of times that the t-statistic was above the critical value of 5%. For the heteroskedastic robust (HR), the normal standard critical values were used. For the SCPC, the critical values were estimated as in Müller and Watson (2022b) using the accompanying post-estimation stata command `scpc`, `cvs`. The HR method uses the `robust` method in Stata. The SCPC uses the post-estimation command `scpc` provided by the authors.

**UPDATE:** This time I randomized the locations. On the previous try, I fixed them. For this iteration, I found a bug in the code on models 2-4, where  $x_l$

Table 1: Rejection frequencies

Model	Replication		M&W (2022)	
	HR	SCPC	HR	SCPC
1	0.50	0.05	0.51	0.05
2	0.48	0.06	0.52	0.08
3	0.47	0.03	0.52	0.15
4	0.33	0.05	0.50	0.08

was constant in all of them. The rejection frequencies went down for almost everyone. I will continue to try what we discussed:

1. Fix locations, and see the results now that the bug is fixed
2. Use intercept in estimation
3. Try different a value for  $\beta$
4. Use `norm` distance instead of `abs`

The rejection frequencies in Table 1 for the SCPC seem to be exactly the expected  $\alpha$  and for model 1 is exactly what Müller and Watson (2022b) show.

### 1.1 Drawing from $\mathcal{G}_{exp}(c)$

To draw from  $\mathcal{G}_{exp}(c_{0.03})$ , first I estimated the matrix of distances  $D$  from the fixed  $s_l$  locations. Then, I found the value of  $c$  such that the average pairwise correlation of the covariance matrix  $\Sigma(c) = \exp(-c * D)$  is close enough to  $\bar{\rho} = 0.03$ . Then, I used the covariance matrix to draw from a multivariate normal distribution with mean zero and covariance matrix  $\Sigma(c)$ .

$$2 \quad e_l \sim \mathcal{G}_{exp}(c_{min})$$

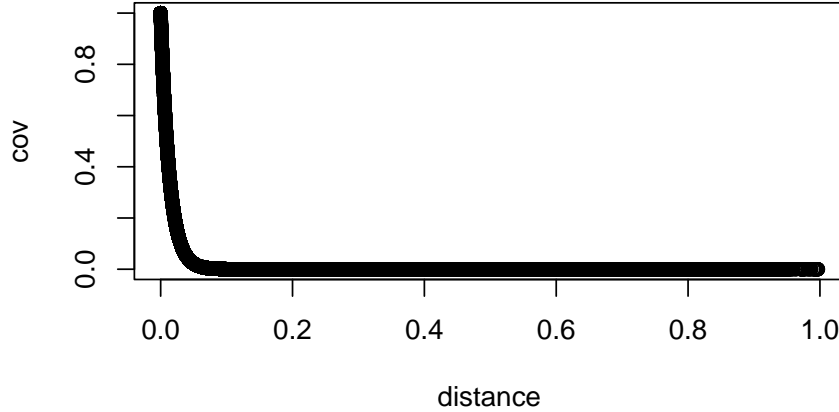


Figure 1: Distance vs Covariance

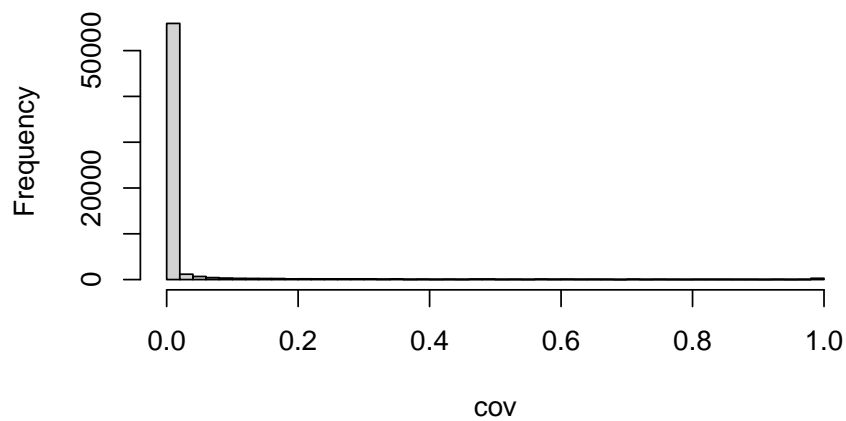


Figure 2: Covariance histogram

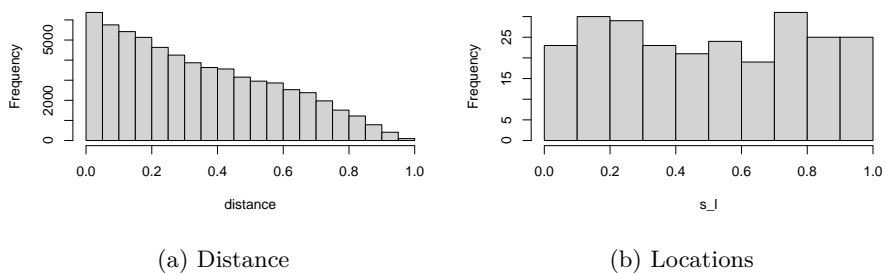


Figure 3: Distance and location histograms

### 3 Power of the test using SCPC

Müller and Watson (2022b) discuss briefly the power of the tests using their variance estimator. They discuss the difference in the power of the tests using the conditional and the unconditional SCPC, the trade-off in the length of the CI and the number of the principal components used to calculate  $\sigma_{SCPC}$ , and the expected length of the CI as a function of  $\bar{\rho}_{max}$ .

#### 3.1 Conditional vs. Unconditional SCPC (p.8)

The authors argue that the size-adjusted power (average length of the confidence intervals) of SCPC and C-SCPC (conditional) are identical because both methods are based on the same  $t$ -statistic and only differ in their critical values. The authors argue that the differences in the average confidence interval lengths of both methods are small. They based their argument on an unreported comparison realized for only one experiment design (p. 8).

The authors acknowledge that the size-adjusted Kernel method is “somewhat” more efficient than SCPC and C-SCPC. The authors argue however that this is not a reason to prefer the Kernel method because in practice the “size-adjustment that adjust for the larger bias in the  $\hat{\sigma}_K^2$ ” is not feasible.

#### 3.2 Trade-off in choosing $q$ (p.4)

The authors discuss that there is a trade-off between the number of principal components  $q$  and the expected length of the 95% confidence intervals. In particular, for a fixed critical value, the expected length of the confidence interval falls as  $q$  increases, but larger  $q$  requires larger critical values to control coverage. The authors suggest that minimizing the expected 95% confidence interval length under the iid benchmark yields a value of  $q$  that works well for a range of values of  $c$ .

#### 3.3 Expected length of confidence intervals as a function of $\bar{\rho}_{max}$ (p.11,12)

The authors also discuss what happens to the rejection frequencies and confidence intervals if the researcher misjudges the spatial correlation  $\bar{\rho}_{max}$ . For example, if the researcher provides a significantly low average pairwise correlation  $\bar{\rho}_{max} = 0.03$  when the true is actually higher  $\bar{\rho} = 0.10$ , the authors argue that the average rejection increases but marginally. The authors performed a set of 14 experiments with different DGPs, first with the true  $\bar{\rho} = 0.03$  and then increased to 0.10, but they kept the provided  $\bar{\rho}_{max} = 0.03$ . They observed that the average rejection frequency increased from 0.04 to 0.07. However, when you look at the supplemental materials the rejection frequency for some of the experiments went up even above 0.4.

The authors also discuss how the confidence intervals increase if the researcher

gives a higher value of  $\bar{\rho}_{max}$  than the true  $\bar{\rho}$ . The authors compare the CI generated by the C-SCPC method to the length of the oracle  $\pm 1.96$  interval length when the  $\bar{\rho}_{max} = 0.1, 0.03, 0.01, 0.003$ . The authors find that the average C-SCPC CI is 1.55, 1.2, 1.1, and 1.05 times larger than the oracle  $\pm 1.96$  CI. The authors conclude that the cost of the default value of  $\bar{\rho} = 0.03$  is about a 20% increase in the CI lengths.

## 4 Summary of the `scpc` stata output

Description of the stata output from the post-estimation command `scpc`:

- Coefficient
- S.E.: The `scpc` post-estimation command displays the standard error as the SCPC estimated variance using eq. 5 modified as follows  $\hat{\sigma}_{SCPC}/(\sqrt{n}\sqrt{q})$  (line 559, `scpc.ado`).
- t: The command reports the t-statistic, using eq. 4, but using the SCPC variance estimator (eq. 5)
- “P>|t|”: This value corresponds to the maximum rejection probability for the adjusted t-statistic. The maximum rejection probability,  $\mathbb{P}_{c \geq c_{max}}(\cdot)$ , is the maximum of the rejection probabilities computed for a range of  $c$  values, that go from  $c_{min}$  to  $c_{max}$ .  $c_{min}$  is provided by the researcher as a maximum average pairwise correlation,  $\bar{\rho}$ . The default is  $\bar{\rho} = 0.03$ .  $c_{max}$  is the maximum value for which size control is checked. In the code,  $c_{max}$  corresponds to a minimum average pairwise correlation of 0.00001 (line 478, `scpc.ado`).
- Confidence interval: The `scpc` command computes the CI using the critical value that corresponds to a 5% maximum rejection probability. For the range of  $c$  values  $\in [c_{min}, c_{max}]$ , the algorithm searches for the critical value that corresponds to a maximum rejection probability  $\alpha = 0.05$ , the desired size of the test.

## 5 Computing rejection probability

$\mathbb{P}_{c \geq c_{max}}(\cdot)$  computes the maximum rejection probability for a range of  $c$  values that go from  $c_{min}$  —which corresponds to  $\hat{\rho}$ — to  $c_{max}$ , a minimum value for which size control is checked. In the code, this value corresponds to a minimum average pairwise correlation of 0.00001 (line 478, `scpc.ado`).

The rejection probability is computed using a Gaussian quadrature and the eigenvalues of  $\Omega(c)$  (pag. 12), for a given critical value.

Then, the critical value is selected such that the rejection probability is the desired, for example,  $\alpha = 0.05$  (line 144, `scpc.ado`)

## Intro

Open questions to address in Spatial Correlation Principal Component (SCPC) (Müller and Watson 2022b, 2022a):

1. What's the “worst case” covariance matrix? How does that work?
2. Can SCPC accommodate a vector of variables? A: Not currently. It only works for scalars. The authors discuss that the main challenge for the multivariate version of the SCPC is specifying the worst-case benchmark model and numerically determining the critical value. In the case of the C-SCPC, the challenge is to specify an appropriate multivariate extension of the heteroscedastic model for  $e_l$  given  $e_l = x_l^s a_l$ .
3. Does the method work if there exists measurement error in the locations?

## 6 The setting

$$y_l = x_l \beta + e_l$$

$\mathbb{E}[e_l|x_l] = 0$ , and  $(y_l, x_l, e_l)$  are associated with observed spatial location  $s_l \in \mathbb{R}^d$ . Spatial location in time is  $d = 1$ , two dimensions, like in altitude and latitude,  $d = 2$ , and so on.

$e_l$  is generated by a Gaussian process with covariance function  $cov(e_l, e_{l'}) = \exp(-c||s_l - s_{l'}||)$ , where  $s_l$  and  $s_{l'}$  denote the spatial locations of  $e_l$  and  $e_{l'}$ , and  $c > 0$  is a parameter that governs the strength of the spatial correlation. The value of  $c$  is calibrated to induce a specific average pairwise correlation  $\bar{\rho} = [n(n-1)]^{-1} \sum_{l, l' \neq l} cov(e_l, e_{l'})$ . In other words,  $c = c_{\bar{\rho}}$  solves  $[n(n-1)]^{-1} \sum_{l, l' \neq l} \exp(-c_{\bar{\rho}}||s_l - s_{l'}||) = \bar{\rho}$ .

Let  $\Sigma(c)$  be the covariance matrix of  $e_l$  evaluated at the sample locations, so that  $\Sigma(c)_{l, l'} = \exp(-c||s_l - s_{l'}||)$ , and let  $\bar{\rho}(c)$  denote the resulting average pairwise correlation  $\bar{\rho}(c) = [n(n-1)]^{-1} \sum_{l, l' \neq l} \Sigma(c)_{l, l'}$ . If the researcher desires a test that controls size for values of  $\bar{\rho}$  as large as  $\bar{\rho}_{max}$ , then he can choose  $c_{min}$  such that  $\bar{\rho}(c_{min}) = \bar{\rho}_{max}$ . Then,  $\Sigma(c_{min})$  is the **worst case** covariance matrix in the sense that it induces the largest value of  $\sigma^2$  among all  $\Sigma(c)$  with  $\bar{\rho} \leq \bar{\rho}_{max}$ .

## 7 Literature Review

### 7.1 Müller and Watson (2022b)

The objective of the paper is to develop a robustified version of their Spatial Correlation Principal Components (SCPC) (Müller and Watson 2022a). In particular, the authors propose modifications to deal with non-stationarities and strong dependence in finite samples (small samples).

According to the authors, the SPCPC method addresses the challenge of spatial correlation robust inference under small samples and *empirically relevant* forms of strong dependence (which ones?).

According to the authors, Conley (1999) doesn't work well in small samples because it relies on the consistency of the estimator of  $\sigma^2$ , while SCPC (and fixed-b type approaches) rely on the *stationarity* of  $u_i$ . Stationarity might break in practice, for example, when  $x_i$  is a dummy for treatment, and treatment is more likely in one region than another region.

The SCPC method is based on a principal component estimator of  $\sigma^2$  based on a pre-specified “worst-case” exponential covariance function conditional on the observed locations.

(On going...)

## References

- Conley, T G. 1999. “GMM Estimation with Cross Sectional Dependence.” *Journal of Econometrics* 92: 1–45.
- Müller, Ulrich K, and Mark W Watson. 2022a. “SPATIAL CORRELATION ROBUST INFERENCE.” *Econometrica* 90: 2901–35. <https://doi.org/10.3982/ECTA19465>.
- . 2022b. “Spatial Correlation Robust Inference in Linear Regression and Panel Models.” *Journal of Business & Economic Statistics* 00: 1–15. <https://doi.org/10.1080/07350015.2022.2127737>.