

# Web-based material for “A climate of uncertainty: accounting for error in climate variables for species distribution models”

By Jakub Stoklosa, Christopher Daly, Scott D. Foster,  
Michael B. Ashcroft and David I. Warton

## Web Appendix A: Algorithms and estimation details

### A1: Monte Carlo expectation–maximization (MCEM)

We used an MCEM algorithm for estimation of errors-in-variables models with classical error. This is a variation on the traditional EM algorithm, useful for fitting models with unobserved components (random effects, latent variables, predictors measured with error, missing data), in which the expectation integral at the E-step is evaluated using Monte Carlo integration. While the MCEM algorithm (as originally proposed [Wei and Tanner \(1990\)](#)) performs a standard Monte Carlo integration sampling from the posterior distribution of prediction error, we used a variation of the method where we sampled from the prior distribution for prediction error and evaluated the integral using importance sampling. This had the advantages that it is simple and stable, given that the Monte Carlo estimates of prediction error do not need to be updated from one iteration to the next, only their importance weights do (which are readily available since they are proportional to the respective marginal likelihood contributions). An advantage of EM-type algorithms is that the maximisation (M-step) usually can be expressed in a simple form, in our case, as a weighted GLM of “complete data” with importance weights. Our algorithm was as follows:

- 1 • We simulate  $B$  replicate Monte Carlo values for prediction error  $U_i$  from the prior
- 2 distribution  $f(X_i)$ , to construct replicates of  $X_i$ , using  $\tilde{X}_i = w_i - U_i$ .
- 3 • The weights are the observations proportional to  $f(y_i | X_i; \beta)f(X_i)$ .
- 4 • Fit a GLM with weights.

5 This procedure is repeated until convergence of  $\beta$  is met. In both the simulations and  
6 case study example, we set  $B = 500$ . See Web Appendix C for further details.

7 Importance sampling works best when the proposal distribution closely matches the  
8 posterior. In our case, there was not a lot of information in the data that could  
9 be used to estimate prediction error, so the prior seemed to work relatively well as  
10 a proposal. This was diagnosed by looking at importance weights, which should be  
11 similar in magnitude across MC samples.

12 Standard errors for  $\hat{\beta}$  are obtained by using the Fisher's information function via  
13 the `vcov()` function in R – i.e. we apply `vcov()` on the final weighted GLM after  
14 convergence is met, see Web Appendix C for further details. We use these to calculate  
15 the standard errors of the linear predictor (i.e.  $\sqrt{X^T \text{Var}(\beta)X}$ ), as plotted in Figure 7.

## 16 **A2: Simulation–extrapolation (SIMEX)**

17 Consider the set of fixed values  $0 = \lambda_1 < \lambda_2 < \dots < \lambda_K$ . We generate new simulated  
18 data by adding independent errors from  $U_i^* \sim N(0, \lambda_k \sigma_u^2)$  for  $k = 1, \dots, K$  to the  
19 explanatory variables  $w_i$ . Note that since we increase  $\lambda_k$  at each succession, we add  
20 more error to the contaminated explanatory variables. For each simulated data set we  
21 fit a GLM-SDM to obtain the estimates for  $\beta$ . We repeat this simulation/estimation  
22 procedure  $B$  times for each succession  $\lambda_k$ , and then calculate the mean values for the  
23 estimates of  $\beta$ . These means are then plotted against the increasing  $\lambda_k$  values, and  
24 polynomial least squares – e.g. a quadratic curve, is fitted to the averaged, error  
25 contaminated estimates. Since the total errors-in-variables variance is  $\sigma_u^2 + \lambda_k \sigma_u^2 =$

1  $(1 + \lambda_k)\sigma_u^2$  for the  $k$ th data set, we extrapolate to the case of no errors-in-variables  
2 (i.e.  $\lambda = -1$ ), which yields the SIMEX estimate.

3 In both the simulations and case study example, we used the R-package `simex` (Lederer  
4 and Kuchenhoff, 2006), and the following default values for  $\lambda = (0.5, 1, 1.5, 2)$  and  $B =$   
5 100, also see Web Appendix C.

## 6 Web Appendix B: Extension to spatial models

### 7 B1: Spatial autocorrelation in the response

8 The joint density in equation (2) assumes the response variable is independent across  
9 observations, and that errors in the explanatory variables are also independent of each  
10 other, leading to “salt-and-pepper” errors for the explanatory variables. To account  
11 for spatial autocorrelation in the response, site-specific random effects are employed  
12 in the linear predictor. Write  $\mu_i = h(X_i^T \beta + \nu_i)$ , where  $\nu_i \sim N(0, \Omega)$  are the random  
13 effects which are independent of  $X_i$ . The covariance matrix  $\Omega = R(\psi)$  models the  
14 spatial autocorrelation, where  $\psi = (\psi_1, \psi_2)$  is a vector of parameters associated with  
15 some fixed matrix  $R$  – e.g.  $R(\psi) = \psi_1 \exp(-d/\psi_2)$ , where  $\psi_1$  is the sill,  $\psi_2$  is the scale  
16 parameter and  $d$  are Euclidean distances, see Li, Tang and Lin (2009) for other spatial  
17 autocorrelation structures. Let  $\theta = (\beta, \psi)$ , then the joint PDF is given by:

$$f(y, w; \theta) = \int \int \left\{ \prod_{i=1}^n f(y_i | X_i, \nu_i; \beta) f(w_i | X_i) f(\nu_i; \Omega) f(X_i) \right\} d\nu dX.$$

18 Both  $\nu$  and  $X$  are treated as missing data and we again consider the MCEM-algorithm.  
19 We followed the same procedure as given in Section A1 but now additionally sim-  
20 ulated replicate Monte Carlo values for the random effects (from the prior distribu-  
21 tion,  $N(0, \Omega)$ ). These simulated random effects were then treated as offsets in the GLM  
22 fit. The parameters  $\psi$  can then be easily estimated outside of the MCEM-algorithm –  
23 e.g. using a profile likelihood approach.

1 In R, the exponential covariance structure to model the spatial autocorrelation, and the  
2 long./lat. values can be used as coordinates within the `rdist.earth()` function (from  
3 the `fields` R-package) to calculate the Euclidean distances. The spatial parameters  $\psi$   
4 can then be estimated using a profile likelihood approach.

## 5 **B2: Spatial errors-in-variables models**

6 The assumption that errors in climate variables are independent is potentially rea-  
7 sonable when the dominant source of errors are within cell variations caused by – e.g.  
8 complex topography or climate-forcing factors operating at a scale finer than the resolu-  
9 tion of the climate grids (Daly, 2006). However, it may be an unreasonable assumption  
10 in relatively flat and homogeneous study areas where within cell variability is low and  
11 errors are dominated by spatially structured errors in how the climate grids are in-  
12 terpolated or predicted. In the case study, we assume an independent structure for  
13 the prediction errors. For completeness, we discuss spatial errors-in-variables models  
14 below, which are technically much more difficult to fit.

15 The only analysis we are aware of which addresses spatially correlated errors-in-variables  
16 models is Li et al. (2009) who dealt with a response variable that followed a Gaussian  
17 distribution. Including spatial components to the prediction error requires an addi-  
18 tional spatial correlation in the true explanatory variable  $X$ . Again we consider a  
19 mixed model approach – i.e. we now write  $X_i = \alpha + \kappa_i + \epsilon_i$ , where  $\kappa_i \sim N(0, \Sigma)$  are the  
20 random effects with some spatial correlation structure  $\Sigma = S(\zeta)$ , and  $\epsilon_i \sim N(0, \sigma_\epsilon^2)$  are  
21 the independent residuals. As in Section B1, the  $\zeta$  is a vector of parameters associated  
22 with some fixed matrix  $S$ , and we must integrate over  $\kappa$  in the joint PDF.

23 Technically this is a more difficult model to implement, and extending this method to  
24 non-normal data, such as using a GLM framework, is a priority for future research.

## 1 **References**

- 2 Daly, C. (2006). Guidelines for assessing the suitability of spatial climate data sets.  
3 *International Journal of Climatology*, 26, 707-721.
- 4 Lederer, W. and Kuchenhoff, H. (2006). A short introduction to the SIMEX and MC-  
5 SIMEX. *R News* **v6.4**, pp. 26–31.
- 6 Li, Y., Tang, H. and Lin, X. (2009). Spatial linear mixed models with covariate mea-  
7 surement errors. *Statistica Sinica* **19**, pp. 1077–1093.
- 8 Wei, G. C. G. and Tanner, M. A. (1990). A Monte Carlo implementation of the EM al-  
9 gorithm and the poor man’s data augmentation algorithms. *Journal of the American*  
10 *Statistical Association* **85**, pp. 699–704.

# 1 Web Appendix C: R-code

```
2 #####
3 #####
4 ## R-code containing functions for the errors-in-variables MCEM method. A
5 ## simulation example is also given, which fit the GLM, SIMEX and MCEM for
6 ## data with errors-in-variables.
7 ##
8 ## Authors: Jakub Stoklosa and David Warton (School of Mathematics and
9 ## Statistics, The University of New South Wales, NSW 2052, Australia).
10 ##
11 ## Please report any problems/suggestions to Jakub Stoklosa at:
12 ## j.stoklosa@unsw.edu.au
13 ##
14 ## This program is meant to be used for non-commercial purposes only.
15 #####
16 #####
17
18 ## Function for implementing the MCEM algorithm for non-Gaussian
19 ## data with errors-in-variables.
20
21 MCEMfit<-function(Y,W,beta,sigma.u,sigma.e,n,B,family,epsilon=0.00001)
22 {
23   set.seed(1000)
24   options(warn=-1);
25   reps<-0;
26   cond<-TRUE;
27
28   ftInit<-glm(Y~W-1,family="binomial");
29   muPred<-rep(predict(ftInit,type="response"),B);
30
31   sigma.u1<-sigma.u[1];
32   sigma.e1<-sigma.e[1];
33
34   U1_j<-rnorm(n*B,0,sd=sqrt(rep(sigma.u1,B)));
35   X1_j<-rep(W[,2],B)-U1_j;
36   X<-cbind(rep(1,B*n),X1_j);
```

```

1
2   while(cond)
3     {
4   ## MC and E-step.
5
6     prX<-dnorm(X1_j,0,sd=sqrt(sigma.e1));
7     bigY<-rep(Y,B);
8
9     prY<-dbinom(bigY,1,muPred);
10
11  ## M-step (updates).
12
13     bigW<-matrix(prY*prX,n,B);
14     sumW<-rep(apply(bigW,1,sum),B);
15
16     mod<-glm(bigY~X-1,weights=as.vector(bigW)/sumW,family="binomial");
17     beta.update<-coef(mod);
18     muPred<-predict(mod,type="response");
19     sigma.e.update<-wt.var(X[,2],w=as.vector(bigW)/sumW);
20
21  ## Convergence monitoring.
22
23     beta.norm<-sum((beta-beta.update)^2);
24     diff.sig_e<-abs(sigma.e.update-sigma.e1);
25
26     reps<-reps+1;   # Keeps track of number of iterations.
27
28     if(diff.sig_e<epsilon && beta.norm<epsilon)
29       {
30         cond<-FALSE;
31         print("convergence met");
32         print(reps);
33         break;
34       }
35
36  ## Update parameters.
37

```

```

1     beta<-beta.update;
2     sigma.e1<-sigma.e.update;
3     }
4
5     beta<-beta.update;
6
7     values<-list(beta=beta,beta.se=sqrt(diag(vcov(mod))));
8     return(values)
9     }
10
11 ## Simulation starts here:
12
13 library("SDMTools");
14 library("simex");
15
16 epsilon<-0.00001; # Convergence threshold.
17 B<-500;          # No. of MC replicates.
18
19 family<-"binomial";
20 n<-500; # Sample size.
21
22 set.seed(1000)
23
24 sigma.e<-1; # True predictor variable variance.
25 x<-rnorm(n,0,sd=sqrt(sigma.e));
26 X<-cbind(rep(1,n),x);
27
28 sigma.u<-0.5; # Variance of uncertainty.
29 w<-x+rnorm(n,0,sd=sqrt(sigma.u)); # Add the uncertainty to explanatory variables.
30 W<-cbind(rep(1,n),w);
31
32 sigma.u<-rep(sigma.u,n);
33 beta<-c(0.5,1);
34
35 mu_Y<-exp(X%*%beta)/(1+exp(X%*%beta));
36 Y<-rbinom(n,1,prob=mu_Y);
37

```



```

1  ## GLM-SDM.
2
3  mod_naiv1<-glm(Y~w,x=TRUE,family=binomial);
4  mod_naiv1;
5  sqrt(diag(vcov(mod_naiv1)));
6
7  # Slope coef: 0.578 (0.087).
8
9  ## MCEM.
10
11 start<-Sys.time();
12 est1<-MCEMfit(Y,W,coef(mod_naiv1),sigma.u,1,n,B,family,epsilon);
13 end<-Sys.time(); end-start;
14 est1;
15
16 # Slope coef: 0.919 (0.117)
17
18 ## SIMEX
19
20 start<-Sys.time();
21 est2<-simex(mod_naiv1,SIMEXvariable="w",measurement.error=sqrt(sigma.u));
22 end<-Sys.time(); end-start;
23 est2;
24 sqrt(diag(est2$variance.jackknife));
25
26 # Slope coef: 0.809 (0.116)
27

```