# GERE

# Some issues in modelling biodiversity using spatially modelled covariates

#### Hideyasu Shimadzu<sup>1</sup>, Scott D. Foster<sup>2</sup>

<sup>1</sup>Marine and Coastal Environment Group, Geoscience Australia. <sup>2</sup>CSIRO Wealth from Oceans Flagship and CSIRO Mathematical and Information Sciences.





Statistical models have enhanced the understanding of the relationship between biodiversity and the environment. Typically, some sort of regression analysis is performed where physical variables are covariates. It is frequently the situation that the covariates are not observed; they are spatial predictions. This study indicates that this process may bias the statistical distribution and the resulting parameter estimates if the variance of the predictions is ignored.



#### Great Barrier Reef data

Data comprises of 1189 sites, where biological and physical variables were measured (Pitcher et al 2007, see Figure 1). The biological outcome variables were presence/absence of a particular species, and species richness (number of species) in a benthic sled sample. The physical variables used in this study were: depth, %carbonate and %mud. We randomly selected 200 sites from this complete set to mimic performing a biological survey (without measuring physical data).



> Figure 1: 200 study sites (orange) randomly selected from 1189 sled survey sites (black).

## Model and approximate bias

Let  $Y_i$  be a biological response such as species presence/ absence or richness, and let  $x_i$  be the vector of physical covariates at the survey site *i*. A generalised linear model is often used with

### Simulation studies

To assess the size of relative bias a simulation study was performed. We simulated biological data given the observed data and analysed them using the spatially predicted covariates  $(\tilde{x}_i)$  and the observed covariates  $x_i$ . A total of 1000 simulations





**Acknowledgements** 

We would like to thank everybody who contributed to the collection of the

Australia's Marine Biodiversity

data. Specific thanks go to Roland Pitcher, Zhi Huang, Kathy Haskard, Mark Palmer, Ross Darnell, Piers Dunstan, Jin Li, and Bill Venables.

#### $\mathrm{E}\left[Y_{i}|\boldsymbol{x}_{i}\right] = h\left(\eta_{i}\right) = h\left(\boldsymbol{x}_{i}^{\top}\boldsymbol{\tau}\right),$

where h is an inverse link function and  $\tau$  is a  $p \times 1$  vector of unknown parameters. However, the physical covariates are commonly not observed but are estimated by spatial predictions,  $\tilde{x}_i$ , based on observations at other sites. Allowing for variance from the spatially predicted covariate, the approximate mean and covariance of the outcomes are

 $\mathrm{E}\left[\mathrm{E}\left[Y_{i}|\tilde{x}_{i}
ight]|\left\{X_{o}
ight\}
ight]=h\left( ilde{x}_{i}^{ op}oldsymbol{ au}
ight)+rac{1}{2}rac{d^{2}h}{dn^{2}}oldsymbol{ au}^{ op}oldsymbol{\Sigma}_{ii}oldsymbol{ au},$ and  $\operatorname{Cov}\left[Y_{i}, Y_{j} | \{\boldsymbol{X}_{o}\}\right] = \operatorname{E}\left[\operatorname{Cov}\left[Y_{i}, Y_{j} | \{\boldsymbol{X}_{o}\}\right]\right] + \frac{dh}{d\eta_{i}} \frac{dh}{d\eta_{i}} \boldsymbol{\tau}^{\top} \boldsymbol{\Sigma}_{ij} \boldsymbol{\tau},$ 

where  $\Sigma_{ii}$  is the cross-covariance of prediction for the *i*th and *j*th physical covariates and  $\{X_o\}$ are the observed physical covariates.

#### were performed and are indexed by k in the equations below.

# Simulation model 1

#### **Simulation model 2** $\mathbf{E}\left[Y_{i}^{(k)}\right] = h\left(\tilde{\boldsymbol{x}}_{i}^{\top}\boldsymbol{\tau}^{(k)}\right); \qquad \mathbf{E}\left[Y_{i}^{(k)}\right] = h\left(\boldsymbol{x}_{i}^{\top}\boldsymbol{\tau}^{(k)}\right).$

#### Results

#### **Relative size of bias**

		Min.	Median	Max.
Presence/absence	<mark>bias</mark> /mean	-19.0%	100.1%	839400.0%
(Cheilostomata Hippaliosina spp)	bias/variance	0.0%	19.2%	1339000.0%
Richness	<mark>bias</mark> /mean	1.3%	3.5%	17.2%
	bias/variance	73.3%	361.6%	1200.0%
Summary statistics of the 200 sites				

### Conclusions

> Using spatially modelled covariates leads to bias in outcome distribution and parameter estimates





> Simulation shows this is a real problem not just a theoretical one

References

> Will need to account for uncertainty in future models

Pitcher R. et al. (2007) Seabed Biodiversity on the Continental Shelf of the Great

Barrier Reef World Heritage Area. CSIRO Marine and Atmospheric Research.

Water, Heritage and the Arts



> Figure 2: Presence/absence – Empirical distribution of estimates from observed physical data (black) and predicted physical data (orange).

the management of Australia's unique environment. (Our stakeholder partners are: AFMA, APPEA, CFA, DAFF, DEWHA, DAFF, the Tourism CRC, and WWF Australia)



**> Figure 3:** Richness – Empirical distribution of estimates from observed physical data (black) and predicted physical data (orange).



**AUSTRALIAN INSTITUTE** OF MARINE SCIENCE