



# Dealing with many correlated covariates in capture–recapture models

Olivier Gimenez<sup>1</sup> · Christophe Barbraud<sup>2</sup>

Received: 28 February 2017 / Accepted: 13 June 2017  
© The Society of Population Ecology and Springer Japan KK 2017

**Abstract** Capture–recapture models for estimating demographic parameters allow covariates to be incorporated to better understand population dynamics. However, high-dimensionality and multicollinearity can hamper estimation and inference. Principal component analysis is incorporated within capture–recapture models and used to reduce the number of predictors into uncorrelated synthetic new variables. Principal components are selected by sequentially assessing their statistical significance. We provide an example on seabird survival to illustrate our approach. Our method requires standard statistical tools, which permits an efficient and easy implementation using standard software.

**Keywords** Animal demography · Population dynamics · Principal-component capture–recapture model · Snow petrel · Survival estimation

## Introduction

Capture–recapture (CR) methods (e.g., Lebreton et al. 1992) are widely used for assessing the effect of explanatory variables on demographic parameters such as survival

(Pollock 2002). Generally however, complex situations arise where multiple covariates are required to capture patterns in survival. In such situations, one usually favors a multiple regression-like CR modeling framework that is however hampered by two issues: first, because it increases the number of parameters to be estimated, incorporating many covariates results in a loss of statistical power and a decrease in the precision of parameter estimates; second, correlation among the set of predictors—aka multicollinearity—may alter interpretation (see below).

To overcome these two issues, Grosbois et al. (2008) recommended to perform a principal component analysis (PCA) on the set of explanatory variables before fitting CR models. PCA is a multivariate technique that explains the variability of a set of variables in terms of a *reduced* set of *uncorrelated* linear combinations of such variables—aka principal components (PCs)—while maximizing the variance (Jolliffe 2002). Grosbois et al. (2008) then expressed survival as a function of the PCs that explained most of the variance in the set of original covariates, typically the first one or the first two ones.

However, the main drawback of this approach is that the PCs are selected based on covariates variation pattern alone, regardless of the response variable, and without guarantee that survival is most related to these PCs (Graham 2003). To deal with this issue in the context of logistic regression, Aguilera et al. (2006) proposed to test the significance of *all* PCs to decide which ones should be retained, instead of a priori relying on the PCs that explain most of the variation in the covariates.

In this paper, we implement the algorithm proposed by Aguilera et al. (2006) to deal with many possibly correlated covariates in CR models, a method we refer to as principal component capture–recapture (P2CR). We apply this new approach to a case study on survival of Snow petrels

**Electronic supplementary material** The online version of this article (doi:10.1007/s10144-017-0586-1) contains supplementary material, which is available to authorized users.

✉ Olivier Gimenez  
olivier.gimenez@cefe.cnrs.fr

<sup>1</sup> CEFE UMR 5175, CNRS, Université de Montpellier, Université Paul-Valéry Montpellier, EPHE, 1919 Route de Mende, 34293 Montpellier Cedex 5, France

<sup>2</sup> CEBC UMR 7372, CNRS-Université de La Rochelle, 79360 Villiers en Bois, France

(*Pagodroma nivea*) that is possibly affected by climatic conditions. In this example, the issue of multicollinearity occurs, and summarizing the set of covariates in a subset of lower dimension is also crucial to get precise survival estimates. Overall, P2CR models can be fitted with statistical programs that perform PCA and CR data analysis. The data and R code are available from GitHub at <https://github.com/oliviergimenez/p2cr>.

## Methods

We used capture–recapture (CR) models to study open populations over  $K$  capture occasions to estimate the probability  $\phi_i$  ( $i = 1, \dots, K-1$ ) that an individual survives to occasion  $i+1$  given that it is alive at time  $i$ , along with the probability  $p_j$  ( $j = 2, \dots, K$ ) that an individual is recaptured at time  $j$ —aka as the Cormack–Jolly–Seber (CJS) model (Lebreton et al. 1992). Covariates were incorporated in survival probabilities using a linear-logistic function:

$$\text{logit}(\phi_i) = \log\left(\frac{\phi_i}{1-\phi_i}\right) = \alpha + \sum_{j=1}^p \beta_j X_{ij} \quad (1)$$

where  $\alpha$  is the intercept parameter,  $X_{ij}$  is the value of covariate  $j$  ( $j = 1, \dots, p$ ) in year  $i$  ( $i = 1, \dots, K-1$ ), and  $\beta_j$  is its associated slope parameter. Covariates were standardized to avoid numerical instabilities. To assess the significance of a covariate in CR models, we used the analysis of deviance (ANODEV; Skalski et al. 1993) that compares the amount of deviance explained by this covariate with the amount of deviance not explained by this covariate, the CR model with fully time-dependent survival serving as a reference. The ANODEV test statistic is given by:

$$\text{ANODEV} = \frac{\text{Dev}(X) - \text{Dev}(\text{constant})}{1} / \frac{\text{Dev}(\text{time}) - \text{Dev}(X)}{K-1} \quad (2)$$

where  $\text{Dev}(\text{constant})$ ,  $\text{Dev}(X)$  and  $\text{Dev}(\text{time})$  stand for the deviance of models with constant, covariate-dependent and time-dependent survival probabilities. To obtain the associated  $P$  value, the value of the ANODEV is compared with the quantile of Fisher–Snedecor distribution with 1 and  $K-1$  degrees of freedom.

To reduce the dimension of the set of covariates ( $X_1, \dots, X_p$ ), we used PCA which aims at finding a small number of linear combinations of the original variables—the principal components (PCs)—while maximizing the variance in ( $X_1, \dots, X_p$ ). Because the variables measurement units often differ, we performed the PCA on the correlation matrix (Jolliffe 2002). To select PCs, we used a forward model selection algorithm as proposed by Aguilera et al. (2006) for the logistic regression. The forward

algorithm begins with no covariates in the model. Each PC is incorporated in simple linear regression-like CR models and the ANODEV  $P$  value calculated. The PC that has the lowest  $P$  value is added to the null model, say  $\text{PC}_k$ . Then the PCs that were not retained are incorporated along with  $\text{PC}_k$  in multiple regression-like CR models, and ANODEV  $P$  values are calculated. In other words, we need to assess the effect of  $\text{PC}_j$  for  $j \neq k$  in the presence of  $\text{PC}_k$  to decide whether  $\text{PC}_j$  should be retained. To do so,  $\text{Dev}(\text{constant})$  and  $\text{Dev}(X)$  are replaced by  $\text{Dev}(\text{PC}_k)$  and  $\text{Dev}(\text{PC}_k + \text{PC}_j)$  in Eq. 2, where  $\text{Dev}(\text{PC}_k + \text{PC}_j)$  is the deviance of the model with survival as a function of both principal components  $\text{PC}_k$  and  $\text{PC}_j$ . We repeat the process until no remaining PC is selected.

All models were fitted using the maximum-likelihood method using MARK (White and Burnham 1999) called with R using package RMark (Laake 2013).

## Case study

The Snow petrel is a medium sized Procellariiform species endemic to Antarctica that breeds in summer. Birds start to occupy breeding sites in early November, laying occurs in early December and chicks fledge in early March. This highly specialized species only forages within the pack-ice on crustaceans and fishes. Data on survival were obtained from a long-term CR study on Ile des Pétrels, Pointe Géologie Archipelago, Terre Adélie, Antarctica. We refer to Barbraud et al. (2000) for more details about data collection. We removed the first capture to limit heterogeneity among individuals, and worked with a total of 604 female capture histories from 1973 to 2002.

The following covariates were included to assess the effect of climatic conditions upon survival variation: sea ice extent (SIE; [http://nsidc.org/data/seaice\\_index/](http://nsidc.org/data/seaice_index/)); air temperature, which was obtained from the Météo France weather station at Dumont d'Urville, as a proxy for sea surface temperature; southern Oscillation Index (SOI) as a proxy for the overall climate condition (<https://crudata.uea.ac.uk/cru/data/soi/>). These environmental variables were averaged over seasonal time periods corresponding to the chick rearing period (January–March: summer period), the non-breeding period (April–June: autumn and July–September: winter), and the laying and incubation period of the same year (October–December: spring). In total, nine covariates were included in the analysis: sea ice extent in summer (SIEsummer), in autumn (SIEautumn), in winter (SIEwinter), in spring (SIEspring), annual SOI, air temperature in summer (Tsummer), in autumn (Tautumn), in winter (Twinter) and in spring (Tspring).

## Results

The CJS model poorly fitted the data ( $\chi^2=221.2$ ,  $df=127$ ,  $P<0.01$ ), and a closer inspection of the results revealed that the lack of fit was explained by a trap-dependence effect (Test2CT,  $\chi^2=103.1$ ,  $df=27$ ,  $P<0.01$ ). Consequently, we estimated two recapture probabilities that differed according to whether or not a recapture occurred the occasion before. By first attempting to simplify the structure of recapture probabilities, we were led to consider an additive effect of time and a trap effect (Electronic Supplementary material, ESM). Estimates of recapture probabilities ranged from 0.14 [standard error (SE) 0.07] to 0.79 (SE 0.09) when no recapture occurred the occasion before and from 0.25 (SE 0.18) to 0.89 (SE 0.09) when a recapture occurred the occasion before (ESM).

Because of multicollinearity, we were led to counter-intuitive estimates of regression parameters in the CR model including all covariates (ESM): the coefficient of SIE in autumn was estimated at 0.5 (SE 0.24) and that of SIE in winter was estimated at  $-0.5$  (SE 0.21) while these two covariates were significantly positively correlated ( $r_p=0.67$ ,  $P<0.01$ ).

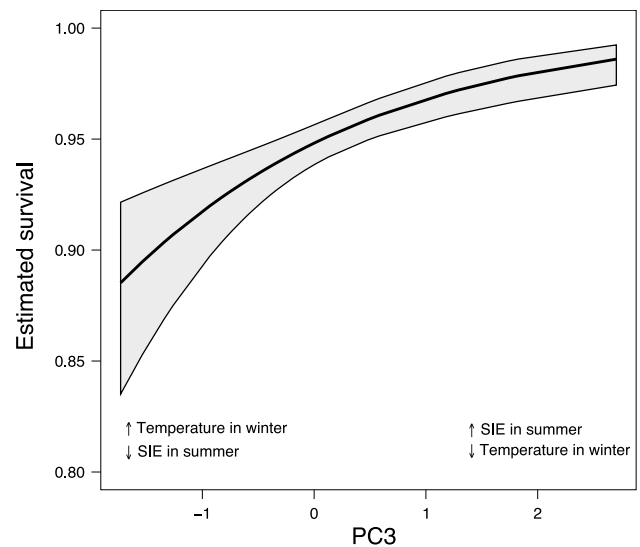
When we applied the P2CR approach, the algorithm selected two PCs, namely PC3 ( $F_{1,27}=7.34$ ,  $P=0.01$ ) at step 1 and PC4 ( $F_{1,26}=4.63$ ,  $P=0.04$ ) at step 2 (ESM), but never did we pick PC1 as we would have done using a classical approach (Grosbois et al. 2008). PC3 was positively correlated to SIE in summer and negatively correlated to temperature in winter, while PC4 was positively correlated to temperature in spring and negatively correlated to SIE in summer (ESM). Survival increased with increasing values of PC3 (Fig. 1), with high values of SIE in summer and low values of temperature in winter (resp. low values of SIE in summer and high values of temperature in winter) corresponding to high (resp. low) survival.

Survival decreased with increasing values of PC4 (Fig. 2), with high values of temperature in spring and low values of SIE in summer (resp. low values of temperature in spring and high values of SIE in summer) corresponding to low (resp. high) survival.

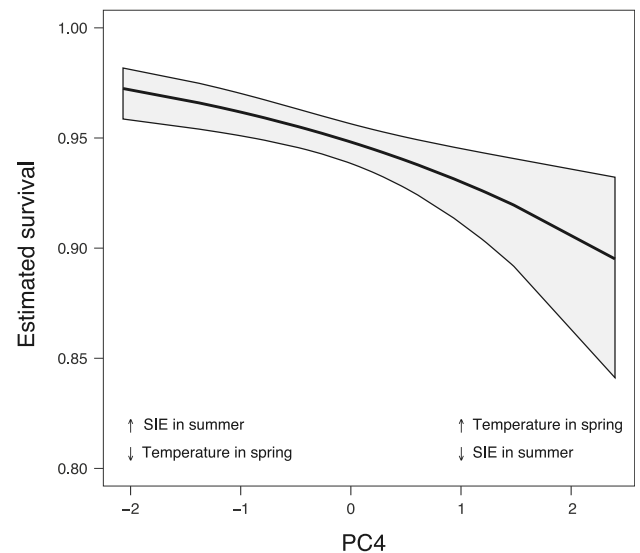
The P2CR approach also led to more precise survival estimates when compared to the model incorporating all original covariates (Fig. 3).

## Discussion

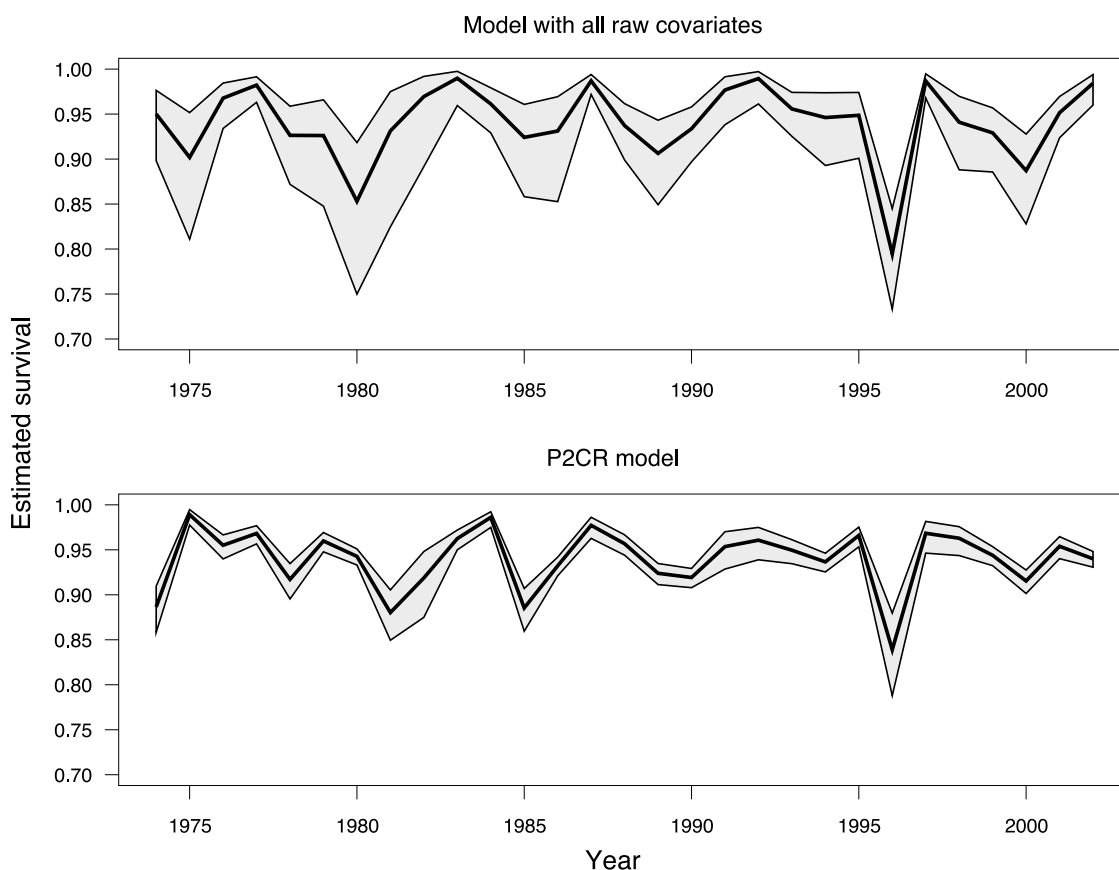
We introduce a new approach combining principal component analysis and capture–recapture models to deal with many possibly correlated explanatory covariates. Our approach requires standard statistical tools, which allows an efficient and easy implementation using standard software.



**Fig. 1** Estimated survival of Snow petrel as a function of PC3 (solid line) with 95% confidence interval (shaded area). Low survival is associated with low values of PC3 that correspond to high values of air temperature in winter and low values of sea ice extent (SIE) in summer; high survival is associated with high values of PC3 that correspond to low values of air temperature in winter and high values of SIE in summer



**Fig. 2** Estimated survival of Snow petrel as a function of PC4 (solid line) with 95% confidence interval (shaded area). High survival is associated with low values of PC4 that correspond to low values of air temperature in spring and high values of sea ice extent (SIE) in summer; low survival is associated with high values of PC4 that correspond to high values of air temperature in spring and low values of SIE in summer



**Fig. 3** Survival of Snow petrel over time as estimated from the model with all original covariates (solid line, top panel) vs. the PC2R model (solid line, bottom panel). 95% confidence intervals are also displayed (shaded area)

### Snow petrels and climatic conditions

In summer, snow petrels exclusively forage within the pack-ice tens to hundreds of kilometers from the colony where they catch sea ice-associated species, such as Antarctic silverfish (*Pleuragramma antarcticum*) and Euphausiids, to feed their chick (Ridoux and Offredo 1989). This is an energetically demanding period for breeding adults and, during years with reduced sea-ice extent, food resources may be less abundant and snow petrels may be forced to cover larger distances to find suitable foraging habitats, with potential survival costs. Assuming air temperature was a proxy of sea surface temperature variations, the negative effect of warmer temperatures on survival is coherent with general patterns found between sea surface temperature and demographic parameters in seabirds (Barbraud et al. 2012). In many marine ecosystems warmer temperatures are associated with decreased primary production and food resources for top predators. Although the low survival in 1996 corresponded to a year with reduced sea-ice extent in summer, the drop in survival was high and remains unexplained at the moment.

### Principal component CR models

When multiple covariates have to be considered to estimate survival, both issues of dimensionality and multicollinearity can lead to biased estimates, inflated precision as well as lack of statistical power. In such a context, the P2CR modeling framework has proved particularly useful in our example, mainly because few PCs were selected which were easily interpretable. We acknowledge that PCs with little interpretability might have been picked up by our method. To make the interpretation easier, PCA results can be post-processed by rotating axes to improve correlations between raw variables and PCs like in the varimax method (Kaiser 1958). Recent developments in the field of multivariate analyses could also be useful, like methods to handle with missing values in PCA (Dray and Josse 2015).

In statistical ecology, one of our objectives is to try and explain variation in state variables such as abundance, survival and the distribution of species. Dimension-reduction methods are promising to deal with many correlated covariates for the analysis of CR or occupancy data.

**Acknowledgements** We greatly acknowledge all of the wintering fieldworkers involved in the monitoring programs in Terre Adélie since 1962, and Dominique Besson for the management of the database. The study on petrels was supported by Expéditions Polaires Françaises, Institut Paul Emile Victor (Program IPEV 109, resp. H. Weimerskirch) and Terres Australes et Antarctiques Françaises.

## References

- Aguilera AM, Escabias M, Valderrama MJ (2006) Using principal components for estimating logistic regression with high-dimensional multicollinear data. *Computational statistics and data Analysis* 50:1905–1924
- Barbraud C, Weimerskirch H, Guinet C, Jouventin P (2000) Effect of sea-ice extent on adult survival of an Antarctic top predator: the snow petrel *Pagodroma nivea*. *Oecologia* 125:483–488
- Barbraud C, Rolland V, Jenouvrier S, Nevoux M, Delord K, Weimerskirch H (2012) Effects of climate change and fisheries bycatch on Southern Ocean seabirds: a review. *Mar Ecol Prog Ser* 454:285–307
- Dray S, Josse J (2015) Principal component analysis with missing values: a comparative survey of methods. *Plant Ecol* 216:657–667
- Graham MH (2003) Confronting multicollinearity in ecological multiple regression. *Ecology* 84:2809–2815
- Grosbois V, Gimenez O, Gaillard JM, Pradel R, Barbraud C, Clobert J, Møller AP, Weimerskirch H (2008) Assessing the impact of climate variation on survival in vertebrate populations. *Biol Rev* 83:357–399
- Jolliffe IT (2002) *Principal component analysis*, 2nd edn. Springer-Verlag, New York
- Kaiser HF (1958) The varimax criterion for analytic rotation in factor analysis. *Psychometrika* 23:187–200
- Laake JL (2013) RMark: An R interface for analysis of capture–recapture data with MARK. AFSC Processed Rep 2013-01, 25 p. Alaska. Fish. Sci. Cent., NOAA, Natl. Mar. Fish. Serv., Seattle
- Lebreton JD, Burnham KP, Clobert J, Anderson DR (1992) Modeling survival and testing biological hypotheses using marked animals: a unified approach with case studies. *Ecol Monogr* 62:67–118
- Pollock KH (2002) The use of auxiliary variables in capture–recapture modelling: an overview. *J Appl Stat* 29:85–102
- Ridoux V, Offredo C (1989) The diets of five summer breeding seabirds in Adélie Land, Antarctica. *Polar Biol* 9:137–145
- Skalski JR, Hoff A, Smith SG (1993) Testing the significance of individual- and cohort-level covariates in animal survival studies. In: Lebreton JD, North PM (eds) *Marked individuals in the study of bird population*. Birkäuser Verlag, Basel, pp 9–28
- White GC, Burnham KP (1999) Program MARK: survival estimation from populations of marked animals. *Bird Study* 46:120–139