

ASSESSING THE POWER TO DETECT SYSTEMATIC CHANGE IN ADÉLIE PENGUIN FORAGING TRIP DURATION

C. Southwell✉, J. Clarke and L.M. Emmerson
Australian Antarctic Division
Department of Environment and Heritage
Channel Highway, Kingston 7050
Tasmania, Australia

Abstract

Models of variability in Adélie penguin foraging trip durations were constructed and fitted to data collected at Béchervaise Island over a 12-year period when only natural variation was known to occur. Variability among trips and penguins was greater in the crèche stage, but variability among years was greater in the guard stage. Estimates of variability were used to explore the power to detect change under particular impact and monitoring scenarios. Power to detect change was greater in the crèche stage than the guard stage. The gain obtained by increasing the number of penguins or trips sampled diminished rapidly when sample sizes were greater than 30 penguins and three trips per penguin. Statistics were developed to test for three forms of change (step, trend and ramp). A test for change based on the difference between pre- and post-impact means generally performed better than a test based on the slope of a trending post-impact change or a joint test of difference and slope. While foraging trip duration is considered to be sensitive to changes in food availability over time scales of days to weeks, because of the high level of natural between-year variation, it would take many years of post-impact monitoring to detect systematic change with high power unless one were willing to relax the Type I error rate to a rate well above the traditional level of 5%. The strategy of including ice cover as a covariate to explain between-year variation in trip duration increased the power to detect change in the guard stage, but the likely dependence between ice cover and fishing activity could confound interpretation and thus, in this case, this strategy is not recommended.

Résumé

Des modèles de variabilité de la durée des sorties alimentaires du manchot Adélie ont été construits et ajustés aux données collectées à l'île Béchervaise pendant une période de 12 années lorsque la variation naturelle était seule en cause. La variabilité entre sorties et entre manchots était plus importante au stade de crèche, mais entre années, elle était plus forte au stade de garde. Des estimations de variabilité sont utilisées pour évaluer la capacité à détecter les changements dans des scénarios particuliers d'impact et de suivi. La puissance de détection des changements s'est révélée plus élevée au stade de crèche qu'au stade de garde. Le gain obtenu en augmentant le nombre de manchots ou de sorties échantillonnées a diminué rapidement lorsque la taille des échantillons était supérieure à 30 individus et trois sorties par manchot. Des statistiques ont été mises au point pour tester trois formes de changement (marche, tendance et rampe). Un test reposant sur la différence entre les moyennes pré et post-impact donne en général de meilleurs résultats qu'un test basé sur la pente d'un changement post-impact en forme de tendance ou un test combinant différence et pente. Alors que la durée de la sortie alimentaire est considérée comme étant sensible aux changements de disponibilité de nourriture sur une échelle temporelle allant de quelques jours à plusieurs semaines, en raison du niveau élevé de variation naturelle d'une année à une autre, il faudrait de nombreuses années de suivi post-impact pour détecter le changement systématique avec une puissance élevée, à moins que l'on soit prêt à adopter, pour le taux d'erreur de Type I, un taux qui soit bien supérieur au taux généralement admis de 5%. La stratégie consistant à inclure la couverture de glace comme covariante pour expliquer la variation entre années de la durée des sorties augmente la puissance de détection des changements au stade de garde, mais la dépendance probable entre la couverture de glace et l'activité de pêche pourrait fausser l'interprétation et c'est pour cette raison que, dans ce cas, cette stratégie n'est pas recommandée.

Резюме

Были построены модели изменчивости продолжительности походов за пищей пингвинов Адели, которые были подобраны к данным, собранным на о-ве Бешервэз за 12-летний период, когда, как известно, наблюдались только естественные изменения. При рассмотрении походов и пингвинов изменчивость была выше на

ясельной стадии, а изменчивость по годам была выше в период присмотра. Оценки изменчивости использовались для изучения возможности выявлять изменение в случае конкретных сценариев воздействия и мониторинга. Вероятность выявления изменений была выше на ясельной стадии по сравнению с периодом присмотра. Преимущество, полученное за счет увеличения числа отобранных пингвинов или походов, быстро уменьшалось, когда размеры выборки превышали 30 пингвинов или три похода на пингвина. Были разработаны статистические показатели для тестирования трех видов изменений (скачкообразное, плавное и медленное). Критерий изменения, основанный на разнице между средними значениями до и после воздействия, в целом дал лучшие результаты, чем критерий, основанный на крутизне тенденции изменения после воздействия, или чем комбинированный критерий разницы и крутизны. Считается, что продолжительность похода за пищей чувствительна к изменениям в наличии пищи в масштабе времени от дней до недель, однако из-за высокого уровня естественной межгодовой изменчивости потребуются многолетние наблюдения для того, чтобы можно было с большой вероятностью выявить систематическое изменение, если не увеличить размер ошибки первого рода до уровня, намного превышающего традиционные 5%. Стратегия включения ледового покрова в качестве ковариаты, объясняющей межгодовую изменчивость в продолжительности походов, повышает вероятность выявления изменений в период присмотра, но возможная зависимость между ледовым покровом и рыбным промыслом может усложнить интерпретацию результатов, и поэтому в данном случае эта стратегия не рекомендуется.

Resumen

Se formularon modelos de la variabilidad de la duración de los viajes de alimentación del pingüino Adelia, que fueron aplicados a los datos recopilados en la Isla Béchervaise durante un período de 12 años, en el cual se sabe que solamente hubo variabilidad natural. La variabilidad entre viajes y entre pingüinos fue mayor durante el período de guardería, pero la variabilidad interanual fue mayor durante el período de cría. Las estimaciones de la variabilidad fueron utilizadas para estudiar la capacidad para detectar cambios bajo ciertas condiciones, en particular suposiciones relativas al impacto y al seguimiento. La capacidad para detectar cambios fue mayor en la etapa de guardería que en la etapa de cría. La ganancia obtenida al aumentar el número de pingüinos o de viajes de alimentación de la muestra disminuyó rápidamente cuando la muestra incluyó más de 30 pingüinos y tres viajes de alimentación por pingüino. Se desarrollaron pruebas estadísticas para detectar tres tipos de cambios (intervalo, tendencia y rampa). Las pruebas para detectar cambios basadas en la diferencia entre los promedios de las variables antes y después del impacto por lo general fueron más efectivas que las pruebas basadas en la pendiente de una tendencia a cambio después del impacto, o una prueba combinada de la diferencia y la pendiente. Si bien se considera que la duración de los viajes de alimentación es muy sensible a los cambios de la disponibilidad diaria y semanal de alimento, debido a la alta variabilidad interanual se tendría que efectuar el seguimiento después del impacto durante muchos años para poder detectar cambios sistemáticos con análisis de alta potencia, a menos que uno estuviese dispuesto a utilizar una tasa de error Tipo 1 mucho menos estricta y bastante mayor que la tasa tradicional de 5%. La estrategia de incluir la cubierta de hielo como covariante para dar cuenta de la variabilidad interanual de la duración de los viajes de alimentación aumentó la capacidad para detectar cambios durante la etapa de cría, pero la posible dependencia entre la cubierta de hielo y las actividades pesqueras podría confundir la interpretación, y por lo tanto no se recomienda tal estrategia en este caso.

Keywords: Adélie penguin, CEMP, change detection, foraging trip duration, power, temporal variability, CCAMLR

Introduction

The time adult penguins spend at sea foraging to provision their chicks is thought to be a sensitive indicator of prey availability over time scales of days to weeks (Cairns, 1987; Croxall et al., 1988). Foraging trip duration is one of the parameters recommended for measurement under

the CCAMLR Ecosystem Monitoring Program (CEMP), which aims to monitor behaviour of krill-dependent predators in order to detect ecosystem changes and differentiate effects of harvesting from those due to natural environmental variability (CCAMLR, 2003). While the likely sensitivity and short response time to changes in prey availability are seen as 'good' characteristics of foraging trip

duration and of ecosystem indicators in general (Landres et al., 1988; Hilti and Merenlender, 2000), such characteristics may also have the disadvantage of showing an inherently high degree of natural variation or 'background noise' that may tend to obscure any signal that monitoring is trying to detect. Examination of the trade-offs between these characteristics is best undertaken through power analyses that take into account the magnitude of all sources of variation affecting the indicator (Southwell et al., 2004).

At the inception of CEMP little was known about the natural level of variability in foraging trip durations of penguins feeding chicks. Initial power analyses (Boveng and Bengtson, 1989; Whitehead, 1989) were necessarily based on a very limited number of years of data, and so could only address the detection of year-to-year differences rather than deviations from normal levels of natural interannual variability (CCAMLR, 1989). However, datasets collected over the past decade or more of CEMP are now sufficiently extensive to enable comprehensive power analyses to be undertaken to explore the sample sizes and the time frames required to detect change of specified size and power.

This manuscript describes the construction of a model of variation in penguin foraging trip duration and the use of a 12-year dataset from Adélie penguins at Béchervaise Island in East Antarctica to estimate the magnitude of the sources of variability in this parameter. These models and variance estimates were then used to estimate the power to detect change under a number of possible impact and monitoring scenarios.

Methods

Scenarios for post-impact change in mean duration of foraging trips

The consequence of an impact on foraging trip time is presumed to result in one of the following changes: (i) a 'step' change in which the mean foraging trip time rises in the first post-impact year and remains constant thereafter, (ii) a 'trend' change in which there is a constant rate of increase across the post-impact years, and (iii) an intermediate 'ramp' change in which there is a constant rate of increase for a number of years, after which the mean level remains constant.

Monitoring scenario

Development of the models below is based on the following monitoring scenario: (i) foraging

trip duration data are collected for each of a consecutive years prior to an impact (pre-impact or 'baseline' data) and b years after that impact (post-impact data) from $2n$ penguins (i.e. n pairs) each making r foraging trips, (ii) the impact may cause a systematic change in mean foraging trip duration in post-impact years, and (iii) the form of post-impact change may be either a step, ramp or trend change.

Foraging trip duration data were obtained over 12 consecutive years (1991/92 to 2002/03) from Adélie penguins breeding at Béchervaise Island (67°35'S 62°49'E), near Mawson Station in East Antarctica. No krill fishery was operating in the region during that time, so these data can be considered to provide pre-impact or baseline data prior to a possible future impact due to a krill fishery.

Modelling variation in foraging trip duration during the pre-impact period

A general model was constructed to include likely and potential sources of variability in individual foraging trip durations under natural pre-impact conditions. Likely sources of variability included years, mating pair membership, gender, penguins and trips.

Foraging trips by Adélie penguins at Béchervaise Island are characterised by many short trips, interspersed with a variable number of longer trips, and the occasional very long trip (Clarke et al., 1998, 2002). Such data are best analysed on a logarithmic scale because it is more useful to express change on a proportional rather than an absolute scale and, for statistical modelling, additivity was found to apply on the logarithmic scale. Analysis of data on the logarithmic scale has been recommended by the CEMP Subgroup on Statistics (SC-CAMLR, 1996).

The general model equation was:

$$y_{ijkm} = \ln f_{ijkm} = M + s_i + p_{ij} + G_k + e_{ijk} + e'_{ijkm}$$

for $i = 1, 2, \dots, a; j = 1, 2, \dots, n_i;$
 $k = 1, 2; m = 1, 2, \dots, r_{ijk}$ (1)

where: f_{ijkm} is the duration of the m th foraging trip for the k th member ($k = 1$ for female, 2 for male) of the j th pair in the i th year, and y_{ijkm} is the natural logarithm of f_{ijkm} ; M is the average duration across all years for all penguins; s_i is a measure of the natural year-to-year variation that was presumed to arise in a random manner from year-to-year; p_{ij} allows for the natural variation among mating pairs in the duration of time they spend foraging; G_k is a gender effect, i.e. $G_1 - G_2$ is the difference in

average trip duration between females and males; e_{ijk} allows for the natural penguin-to-penguin variation in trip duration within a mating pair that is not due to gender; and e_{ijkm} allows for natural variation in duration among successive trips by the same penguin. The components s_i , p_{ij} , e_{ijk} and e_{ijkm} are presumed to be random values from normal distributions with mean 0 and variances σ_s^2 , σ_{pr}^2 , σ_p^2 and σ_i^2 respectively. The variables are assumed to be independent except for the set of trip components for an individual penguin, the e_{ijkm} terms, which might be expected to be correlated.

It is noted that the model was augmented when initially applied to the available pre-impact data by including a systematic component representing a possible trend over years in order to check for evidence of a systematic change in mean foraging trip duration in the pre-impact period. There was no evidence for a trend. Hence variance estimation was based on equation (1).

This general model was assumed to be the simplest representation of variation in foraging trip duration in the guard stage, when pair members are alternately foraging and guarding chicks. In the crèche stage, however, when both parents are able to forage at the same time, the connection in trip duration between pair members may no longer be present, and a simpler model was applied, namely,

$$y_{ikm} = \ln f_{ikm} = M + s_i + G_k + e_{ik} + e_{ikm}$$

for $i = 1, 2, \dots, a$; $k = 1, 2, \dots, 2n_i$; $m = 1, 2, \dots, r_{ik}$ (2)

where the components have the same meaning as in equation (1). Modelling of the data showed that the assumption of independence between pair members in trip durations during the crèche stage was supported; hence equations (1) and (2) were used to estimate variance components for the guard and crèche stages respectively.

One additional potential source of variability is the extent of ice cover offshore from the breeding colony, which may influence foraging trip duration through its effect on prey availability and/or rate of travel (Irvine et al., 2000; Clarke et al., 2002). This source of variation can be represented by replacing the random year component s_i in equations (1) and (2) as follows:

$$s_i = \gamma(x_i - \bar{x}) + s_i', \quad (3)$$

where x_i is the percentage ice cover in the i th year, \bar{x} is the average of the percentage ice cover across the years for which data were available, and s_i' is

the remaining (unexplained) natural year-to-year variation. The model assumes a linear relation between mean trip duration per year and percentage ice cover.

The primary set of statistics for yearly comparisons is the set of yearly means for foraging trip duration $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_{a+b}$. The \bar{y}_i terms are normally distributed with variance represented by equation (4) for the guard stage and equation (5) for the crèche stage:

$$\sigma_{\bar{y}}^2 = \sigma_s^2 + \frac{1}{n}\sigma_{pr}^2 + \frac{1}{2n}\sigma_p^2 + \frac{\sigma_i^2}{2nr}, \quad (4)$$

$$\sigma_{\bar{y}}^2 = \sigma_s^2 + \frac{1}{2n}\sigma_p^2 + \frac{\sigma_i^2}{2nr}. \quad (5)$$

Estimating variation in foraging trip duration during the pre-impact period

Model equations (1) and (2) were fitted to the Adélie penguin trip data for guard and crèche trips respectively to obtain estimates of the variance components in the pre-impact period. The estimates of variance components were then inserted into equations (4) and (5) to provide an estimated variance for yearly means. As preliminary analysis showed no evidence of a systematic trend over the pre-impact years, all year-to-year variation was assumed to be random.

Modelling variation in foraging trip duration during a post-impact period

Development of models for the post-impact period was based on the following assumptions: (i) equations (1) and (2) represent the contributions to trip duration from various sources of variation for the guard and crèche stages respectively, (ii) an impact after a years causes either a step, ramp or trend change in foraging trip duration across the b post-impact years, but there is no change in variability, and (iii) the nature of the impact is such that any change is uni-directional and positive. With respect to this last point, one would expect a decline in prey availability to affect foraging trip duration in this manner. The effects of step, ramp and trend changes on power were explored because it is not presently obvious how quickly a future krill fishery might impact on ecosystem components. These three forms were chosen to represent two extremes and an intermediate rate of impact. The assumption of a constant variance (on the logarithmic scale) between pre- and post-impact years is the simplest of numerous possible scenarios. Exploration

of other variance scenarios could be undertaken in this modelling framework but is beyond the intent and scope of this paper.

Step change model

In this model it is assumed that mean foraging trip duration is constant for years $i = 1, 2, \dots, a$ prior to an impact (i.e. $M_i = M_0$), then in the year following an impact mean trip duration may increase by a proportion p and remain at that changed level thereafter (i.e. $M_i = M_0(1 + p)$ for $i = a+1, a+2, \dots, a+b$ thereafter). This can be equivalently expressed on a logarithmic scale as:

$$\begin{aligned} \log M_i &= \log M_0 & \text{for } i = 1, 2, \dots, a \\ &= \log(M_0) + \delta & \text{for } i = a+1, a+2, \dots, a+b \end{aligned}$$

where δ is the difference between pre- and post-impact periods on the logarithmic scale (i.e. $\delta = \log(1 + p)$). An estimate of δ , denoted by $\hat{\delta}$, can be obtained as the difference from analysis of data on the logarithmic scale, in the mean for years $a+1$ to $a+b$ minus the mean for years 1 to a , and an estimate of the percentage change in mean trip duration as $100[\exp(\hat{\delta}) - 1]$.

Trend change model

This model assumes a proportional increase in foraging trip duration in the years after an impact (i.e. $M_{i+1} = M_i + pM_i$, or $\frac{M_{i+1}}{M_i} = 1 + p$, where M_i is the expected trip duration in the i th year and p is the proportional increase). As equations (1) and (2) model the logarithm of mean foraging trip duration, the comparison of the mean trip duration in successive post-impact years could be expressed as $\log(M_{i+1}) - \log(M_i)$ or $\log\left(\frac{M_{i+1}}{M_i}\right)$. The requirement for the mean foraging trip duration in each post-impact year to increase by a proportion p over the previous year can therefore be expressed as $\log(M_{i+1}) - \log(M_i) = \log(1 + p)$, or, by assuming a linear regression relation based on the logarithms of the mean trip durations for the years, as $\log(M_i) = \alpha + \beta x_i$. In this representation $\beta = \log(1 + p)$ and the x_i term represents years. Under this model, evidence for the presence of a trend change can be derived from a test of the hypothesis $\beta = 0$ and an estimate of the percentage change determined as $\hat{p} = 100 \exp(\hat{\beta})$, where $\hat{\beta}$ is the estimate of the regression coefficient β .

Ramp change model

This model assumes a proportional increase in foraging trip duration for the first c post-impact years followed by a constant level for the remaining $b-c$ years (i.e. $M_{i+1} = M_i + pM_i$ for $i = a+1, a+2, \dots, a+c$, and $M_i = M_c = M_0(1+p)^c$ for $i = a+c+1, a+c+2, \dots, a+b$, where M_0 is the pre-impact mean level).

Power to detect a change between pre- and post-impact data

Tests for the detection of step and trend changes that are based on tests of δ and β , and whose powers can be examined analytically, are presented below.

'Difference' test

Under the assumption of a step change, if \bar{y}_{pre} and \bar{y}_{post} are the means of the pre-impact and post-impact year means, then $\bar{y}_{post} - \bar{y}_{pre}$ is normally distributed with mean δ and variance

$$\sigma_a^2 = \sigma_y^2 \left(\frac{1}{a} + \frac{1}{b} \right). \quad (6)$$

Evidence for the presence of step change can be derived from a test of the null hypothesis $\delta = 0$, with the test based on the statistic $t = (\bar{y}_{post} - \bar{y}_{pre}) / s_a$, where s_a^2 is an estimator of σ_a^2 defined in equation (6). A value for s_a^2 was obtained by inserting estimated values for the variance components based on pre-impact data into equations (4) and (5). Under the hypothesis of no change up to year a (i.e. the pre-impact years) versus the alternative of an increase in mean of $P\%$ in year $a+1$ (i.e. the first post-impact year) and a constant level thereafter, the statistic t has a distribution under the null hypothesis that is well approximated by a t -distribution with $a-1$ degrees of freedom.

'Slope' test

Evidence for the presence of trend change can be derived from a test of the null hypothesis $\beta = 0$, with the test based on the statistic $t = \hat{\beta} / s_{\hat{\beta}}$, where $\hat{\beta}$ is the standard linear regression estimator of β from the regression of mean log trip duration values for year on a variable representing years. The statistic $s_{\hat{\beta}}^2$ is an estimator of $\sigma_{\hat{\beta}}^2 = \sigma_y^2 / \sum_{i=a}^{a+b} (x_i - \bar{x})^2$. Estimates of the variance components required in equations (4) and (5) were obtained from the analysis based on available pre-impact data.

Under the hypothesis of no change between pre- and post-impact years versus the alternative of

a constant, increasing rate of change of $P\%$ (i.e. the proportionate increase, $p = P/100$) after an impact, the statistic t has a distribution well approximated by a t -distribution with $a-1$ degrees of freedom, where a is the number of pre-impact years. If the minimum percentage increase from one year to the next that is considered to be of significance is set at P_0 (which implies $\beta_0 = \log(1+P_0/100)$), and the chance that the test will incorrectly claim there is a difference (a Type I error) is set at α , then the power of the test is $\Pr(t(a-1, \lambda) > t_{\alpha}(a-1))$ where $\lambda = \beta_0 / s_{\beta}$ (the non-centrality parameter), $t(a-1, \lambda)$ is a non-central t distribution with degrees of freedom $a-1$ and non-centrality parameter λ , and $t_{\alpha}(a-1)$ is the value for a t -distribution with $a-1$ degrees of freedom that is exceeded with a probability α if interest lies only in detecting an increasing trend in mean foraging trip duration across years.

It is noted that even under a trend change scenario a test based on the slope statistic is not necessarily more powerful than the test based on the difference statistic. The reason lies in the fact that the difference statistic makes use of all the pre-impact data through the baseline mean, whereas the slope statistic only makes use of the last pre-impact datum as one of the values for estimation of post-impact slope.

A test for detecting a ramp change

Under the assumption of a ramp change scenario, in practice the number of years for the maximum level to be reached (c) is unknown. Thus there are two unknown parameters: the rate of increase in mean trip duration p , and the number of years before the maximum mean trip duration is reached c (or equivalently the maximum mean trip duration M_c). In this situation there is no optimal test statistic to employ, given that a test would be required for each year from the first post-impact year.

A possible approach is to employ a test that jointly uses the difference and slope tests defined above. In the period of increasing means it might be expected that the slope statistic would be more likely to detect change, and thereafter the difference statistic would be more likely to detect change. This involves computing both the difference and slope statistics for step and trend changes and declaring evidence of change if at least one of these is significant. Unlike the tests for slope and difference individually where power is determined analytically, the power of the joint test can only be determined by simulation. In addition, determining a Type I

error rate of α for a joint test requires that the individual tests are assigned critical values that are less than α .

Power analysis scenarios

Power analyses were conducted for a range of possible impact and monitoring scenarios.

With regard to impact scenarios, an effect size of a 25% increase from the pre-impact baseline mean was chosen as a level that may feasibly occur and is likely to be biologically important. Foraging trip durations have been recorded at this level during some years at Béchervaise Island, and were associated with significant breeding failure (Clarke et al., 2002). Given this effect size, power analyses were undertaken when (i) the effect size is reached immediately after an impact and thereafter remains constant (step change), (ii) the effect size is reached after 10 years of a gradual linear increase (2.26% increase per year), and thereafter remains constant (ramp change), and (iii) the effect size is reached after 30 years (0.75% increase per year; trend change), to simulate scenarios where an impact, such as fishing, may occur suddenly at high levels or trend up to high levels at different rates (analogous to pulse and press impacts respectively as described in the environmental literature (Bender et al., 1984; Underwood, 1991)).

With regard to monitoring scenarios, analyses were undertaken for a varying number of years of post-impact monitoring (1–10, 15, 20, 25, 30), and with and without ice cover as a covariate. Power calculations were undertaken using the difference, slope and joint tests for α levels of 0.05, 0.10 and 0.20, and were based on one-tailed tests. Each of the three tests was applied for each of the three types of change scenarios, because in the real monitoring situation the form of change would not be known in advance of an impact and may not become apparent for many years after the impact, but there would be a need to test for change using a specific statistic immediately after the impact occurs. To achieve these Type I error rates for the joint test it was necessary to set lower levels of α for the individual tests, given that a significant result occurs if at least one of these tests was significant. It was established through trial and error that α had to be set at 0.03, 0.06 and 0.12 for individual tests to obtain, for the joint test, a Type I probability of approximately 0.05, 0.10 and 0.20 respectively that at least one test would provide a significant result when there is no change between pre- and post-impact years.

Note that power analysis for the joint test was undertaken using simulation, whereas analytic calculation was employed for power analysis based on the individual tests. The simulation was undertaken using the R statistical computing package with each estimate of power based on 5 000 simulations.

Data collection

Foraging trip duration data were collected by means of an automated penguin monitoring system (APMS) which automatically recorded the time that uniquely tagged individuals departed from, and arrived at, their breeding colony when undertaking foraging trips. The APMS and tag-implantation methods are described in detail in Kerry et al. (1993a) and Clarke and Kerry (1998). Once individuals had been tagged, trip durations were measured for the same individual over several successive years while it continued to breed in the same colony. Penguins were sexed when originally tagged using the methods described by Sladen (1978). Membership of breeding pairs was determined by scanning incubating penguins at each occupied nest in the colony at two points in time during the breeding period when primarily only males, then females, were present (Kerry et al., 1993b). Trip durations of breeding penguins provisioning chicks were extracted from the APMS dataset as described in Clarke et al. (1998, 2002, 2006). The number of pairs, penguins and trips for which data were obtained across the 12 years ranged from 77–131 (median 95) and 138–247 (median 168) to 384–3437 (median 1715) respectively.

Ice data were derived from the National Snow and Ice Data Center (Comiso, 1990, updated current year) found at www.nsidc.org/data/nsidc-0002.html, and collated as percentage ice cover in a 100 x 100 km square north of the colony averaged over all days in January each year.

Results

Variance estimation

Estimated variance components for year, pair, penguin and trip are shown in Table 1. Trip duration showed the greatest variability, followed by penguin, pair and year. Variability among trips and penguins was greater in the crèche stage than in the guard stage, but variability among years was greater in the guard stage.

Effect on variance of varying the number of penguins and trips

In a monitoring program, the primary variables under the control of the investigator are the number of penguins and trips sampled. Increasing the numbers of penguins and trips will improve the precision of yearly mean estimates, but it is the between-year variance component that is most important and that ultimately determines the ability to detect any systematic change from natural variation. Table 2 shows how the precision of yearly mean trip duration estimates is improved (i.e. variance is reduced) as the numbers of penguins and/or trips per penguin are increased. The patterns are generally similar for the guard and crèche stages. The inclusion of penguins that make only one or two trips substantially increases the variance of the yearly mean, but using more than three trips per penguin only marginally reduces the variance. Given data from three trips per penguin, increasing the number of penguins from 30 to 50 results in a small reduction in variance, but a further increase in sample size beyond 50 penguins returns only very marginal improvement. All results reported below are for sample sizes of 30 penguins and three trips per penguin.

Power in relation to impact and monitoring scenarios

Figure 1 compares the power of tests based on the difference, slope and joint statistics to detect step, ramp and trend changes in the guard and crèche stages. There was generally close agreement, for all types of change, between the power of tests based on the difference statistic and the joint test based on both statistics, although the test based on difference was slightly more powerful than the joint test with a step change. It might be expected that the joint test, which utilises both the difference statistic and the slope statistic, would strengthen the detection of a trend while continuing to have power to detect a step change. The reason this is not the case lies in the fact that the Type I error rate for each individual test must be reduced to ensure that the nominated Type I error rate is maintained for the joint test. Hence there is a lower chance for the individual statistics to detect change in the area in which they are superior. There is no scenario under which the slope test proves substantially superior, and for the trend and ramp scenarios it has much lower power than the other tests.

Power to detect an increase in foraging trip duration using only the difference statistic for selected impact and monitoring scenarios is shown in Figure 2. The most obvious finding is that a

gradual trending increase of 0.75% per year over 30 years (i.e. 25% increase after 30 years) is, in practical terms, almost impossible to detect in that time frame, with the probability of detection exceeding 0.80 only in the single case of 30 years of monitoring during the crèche period with a Type I error rate of 0.20. A step increase of 25% is, as one would expect, much easier to detect, but even in this case detecting a change is not always easy and depends largely on the Type I error rate one is willing to accept; for example setting the chance of a Type I error to 0.05 means that it would take 10–15 years to detect a change with probability >0.80 for the crèche and guard stages respectively, whereas if one were willing to relax the Type I error rate to 0.20, changes could be detected with probability 0.80 in three to four years. If an impact caused a ramp increase of the magnitude and rate investigated and the difference statistic was used to test for a change, one would have to accept a Type I error rate of 0.20 and monitor for up to 10–15 years after the impact to detect a change with >0.80 probability.

Ice cover as a covariate

The mean percentage ice cover in a 100 x 100 km square offshore from the colony during January was an important explanatory variable for the guard stage but not for the crèche stage. The year variance component during the guard stage was reduced by approximately 30% when percentage ice cover was fitted as a linear predictor. This improvement in the capability of tests to detect systematic change is illustrated in Figure 3, which shows that a power of 0.80 is reached after 10 years of post-impact monitoring when ice is considered as a covariate compared with 15 years when ice is not considered.

Discussion

One of the difficulties in a real monitoring situation is that the form of change that may occur due to an impact is unknown prior to the impact occurring, and may not become apparent until many years after the impact has occurred. This can present difficulties in developing and applying an appropriate or optimal test. The issue of uncertainty in the form of change to be tested has not been well addressed in the environmental monitoring literature. This study indicates that, of the difference, slope and joint tests investigated, the difference test performs best over a wide range of scenarios. A recommendation from this study is that the difference test be used in preference to slope and joint tests in any future testing of change.

An important result that emerges from the modelling and power analysis of Adélie penguin foraging trip duration data is that the high level of between-year variability, naturally present for this parameter at Béchervaise Island, makes detection of systematic change over a short period of time very unlikely unless one is willing to relax the Type I error rate to well above the traditional level of 5%. Indeed, if prevention of change is to be achieved within the time frame of two to three decades specified in Article 2 of the CAMLR Convention (CCAMLR, 2004), then there may be no option but to accept Type I error rates substantially greater than traditional levels. This would be consistent with an increased questioning of the use of the 5% significance level for null hypothesis tests in environmental monitoring (e.g. Millard, 1987; Fairweather, 1991; Peterson, 1993; Skalski, 1995). Implicit in the traditional convention is the assumption that Type I errors are more important than Type II errors, whereas many authors argue that in the field of environmental monitoring and management the opposite is in fact the case, because the failure to detect a real change may have catastrophic environmental and remedial costs that far outweigh the cost of investigating an occasional false alarm (Peterson, 1993; Gibbs et al., 1999). Millard (1987) and Keough and Mapstone (1997) recommend an alternative approach in which the costs (financial, environmental and/or social) of making Type I and Type II errors are quantified and the levels for each type of error are balanced to reflect these costs. An important outcome of this study is that it focuses attention on the interaction and balance between Type I and Type II errors. Making a decision on where the balance lies however, is a policy exercise rather than a statistical one.

Associated with the finding of high between-year variation is the consequence that a strategy of increasing the number of penguins and/or trips sampled each year to improve power is increasingly ineffective beyond reasonably low levels. For Adélie penguins at Béchervaise Island, the gain rapidly diminishes beyond sample sizes of 30 penguins and three trips per penguin. The original (still current) recommended sampling regime for CEMP parameter A5 (duration of foraging trips) is to calculate bird means of multiple trips throughout chick rearing from 20 or more penguin pairs (CCAMLR, 2003). This sampling regime makes no recommendation on the number of trips per bird. The results presented here indicate that this original recommendation, which was based on very limited knowledge of between-year variability, is reasonable but somewhat low for Adélie penguins at Béchervaise Island. Analysis of other long-term datasets that have resulted from the application

of CEMP over the past one to two decades would allow an assessment of whether the specific results reported here apply more generally to other sites and species.

Sample sizes employed at Béchervaise Island for the monitoring of foraging trip duration since CEMP commenced there in 1990/91 have been much higher than 30 penguins and/or three trips per penguin (see 'Methods'). While these large sample sizes have been useful for gaining an understanding of the ecology of foraging behaviour, they would not be effective in improving the power to detect a future systematic change due to a krill fishery or any other impact. It would be possible to reduce the number of penguins sampled in future monitoring without any deleterious effect on the power to detect a change, provided a minimum of three trips is measured for each bird.

It was hoped that a strategy of including covariates in the statistical model might increase the power to detect change by explaining a proportion of the between-year variability, and this was confirmed for the covariate ice cover in the guard stage. However, this strategy is not recommended, despite some gain in statistical power, because the covariate of ice cover is unlikely to be independent of the change that monitoring is trying to detect. For example, decreasing ice cover may increase the accessibility of locations to fishing, or may impact directly on krill availability by decreasing the extent of krill habitat. If these relationships exist, adjusting foraging trip duration for changes in ice cover may hinder interpretation by camouflaging detection of real change in foraging trip duration.

The results of this study emphasise the trade-off that exists in using a highly 'sensitive' indicator to detect change if that indicator, by the nature of its sensitivity, has a high degree of inherent natural variation. Thus, while foraging trip duration is considered sensitive to changes in food availability on time scales of days to weeks (Croxall et al., 1988; CCAMLR, 2003), this trait may have so much associated natural variation that it may take several years to confidently distinguish a systematic change from the noise of natural variability.

Although the power analyses carried out in this study pertain to a single index within a multi-parameter monitoring program at Béchervaise Island, the general findings are relevant to any single or combined index that shows a high degree of inherent natural variation. Understanding the variability and power associated with each component of a combined index is important in the context of

detection of change, both from an interpretive point of view and as an adjunct to managing ambiguities resulting from missing data.

Conclusions

The high level of between-year variability in foraging trip duration that is naturally present at Béchervaise Island limits the power to detect systematic change over short time periods unless one is willing to relax the Type I error rate to well above the traditional level of 5%. Future population modelling is required to ascertain the demographic effects of detectable increases in foraging trip duration, given reasonable correlations with fecundity and/or survival. This will determine whether such effect sizes are of sufficient magnitude to be useful in the context of Article 2 of the CAMLR Convention, i.e. prevention of changes or minimisation of the risk of changes in the marine ecosystem which are not potentially reversible over two or three decades (CCAMLR, 2004). If significant change is not detectable within this time frame, then risk analysis rather than data-driven processes will be required to enable decisions to be made at the management level.

Acknowledgements

The authors are grateful for the assistance of many field staff and engineers, particularly Lyn Irvine, Megan Tierney and Kym Newbery. Thanks are extended also to Glen McPherson, who developed the statistical models and carried out the power analyses, and to Ben Raymond, who collated the sea-ice data. Comments by Glen McPherson, Keith Reid and an anonymous reviewer improved the manuscript.

References

- Bender, E.A., T.J. Case and M.E. Gilpin. 1984. Perturbation experiments in community ecology: theory and practice. *Ecology*, 65 (1): 1–13.
- Boveng, P. and J.L. Bengtson. 1989. On the power to detect changes using the standard methods for monitoring parameters of predatory species. In: *Selected Scientific Papers, 1989 (SC-CAMLR-SSP/6)*. CCAMLR, Hobart, Australia: 377–397.
- Cairns, D.K. 1987. Seabirds as indicators of marine food supplies. *Biol. Oceanogr.*, 5: 261–271.

- CCAMLR. 2004. *CCAMLR Basic Documents*. CCAMLR, Hobart, Australia: www.ccamlr.org/pu/e/e_pubs/bd/toc.htm.
- CCAMLR. 1989. Instructions for the preparation of sensitivity analyses. Document *WG-CEMP-89/13*. CCAMLR, Hobart, Australia. (Prepared by the Secretariat and the Convener of the Working Group on CEMP.)
- CCAMLR. 2003. *CCAMLR Ecosystem Monitoring Program: Standard Methods for Monitoring Studies*. CCAMLR, Hobart, Australia.
- Clarke, J. and K. Kerry. 1998. Implanted transponders in penguins: implantation, reliability, and long term effects. *J. Field Ornithol.*, 69: 149–159.
- Clarke, J., B. Manly, K. Kerry, H. Gardner, E. Franchi, S. Corsolini and S. Focardi. 1998. Sex differences in Adélie penguin foraging strategies. *Polar Biol.*, 20 (4): 248–258.
- Clarke, J., K.R. Kerry, L. Irvine and B. Phillips. 2002. Chick provisioning and breeding success of Adélie penguins at Béchervaise Island over eight successive seasons. *Polar Biol.*, 25 (1): 21–30.
- Clarke, J., L. Emmerson and P. Othahal. 2006. Environmental conditions and life-history constraints determine foraging range in breeding Adélie penguins. *Mar. Ecol. Progr. Ser.*, 310: 247–261.
- Comiso, J. 1990 (updated current year). *DMSP SSM/I daily polar gridded sea ice concentrations*, June to September 2001. Maslanik, J. and J. Stroeve (Eds). Boulder, CO: National Snow and Ice Data Center. Digital media.
- Croxall, J.P., T.S. McCann, P.A. Prince and P. Rothery. 1988. Reproductive performance of seabirds and seals at South Georgia and Signy Island, South Orkney Islands 1976–1986: implications for Southern Ocean monitoring studies. In: Sahrhage, D. (Ed.). *Antarctic Ocean and Resources Variability*. Springer-Verlag, Berlin Heidelberg: 261–285.
- Fairweather, P.G. 1991. Statistical power and design requirements for environmental monitoring. *Aust. J. Mar. Fresh. Res.*, 42: 555–567.
- Gibbs, J.P., H.L. Snell and C.E. Causton. 1999. Effective monitoring for adaptive wildlife management: lessons from the Galapagos Islands. *J. Wildl. Mgmt.*, 63 (4): 1055–1065.
- Hilti, J. and A. Merenlender. 2000. Faunal indicator taxa selection for monitoring ecosystem health. *Biol. Cons.*, 92 (2): 185–197.
- Irvine, L.G., J.R. Clarke and K.R. Kerry. 2000. Low breeding success of the Adélie penguin at Béchervaise Island in the 1998/99 season. *CCAMLR Science*, 7: 151–167.
- Keough, M.J. and B.D. Mapstone. 1997. Designing environmental monitoring for pulp mills in Australia. *Water Science Technology*, 35 (2–3): 397–404.
- Kerry, K.R., J.R. Clarke and G.D. Else. 1993a. The use of an automatic weighing and recording system for the study of the biology of Adélie penguins (*Pygoscelis adeliae*). *Proc. NIPR Symp. Polar Biol.*, 6: 62–75.
- Kerry, K., J. Clarke and G. Else. 1993b. Identification of sex of Adélie penguins from observation of incubating birds. *Wildl. Res.*, 29: 725–732.
- Landres, P.B., J. Verner and J.W. Thomas. 1988. Ecological uses of vertebrate indicator species: a critique. *Cons. Biol.*, 2 (4): 316–328.
- Millard, S.P. 1987. Environmental monitoring, statistics, and the law: room for improvement. *American Statistical Association*, 41: 249–253.
- Peterson, C.H. 1993. Improvement of environmental impact analysis by application of principles derived from manipulative ecology: lessons from coastal marine case histories. *Aust. J. Ecol.*, 18: 21–52.
- SC-CAMLR. 1996. Report of the Subgroup on Statistics. In: *Report of the Fifteenth Meeting of the Scientific Committee (SC-CAMLR-XV)*, Annex 4, Appendix H. CCAMLR, Hobart, Australia: 251–273.
- Skalski, J.R. 1995. Statistical considerations in the design and analysis of environmental damage assessment studies. *J. Envir. Manage.*, 43: 67–85.
- Sladen, W.J.L. 1978. Sexing penguins by cloacoscope. *International Zoo Yearbook*, 18: 77–80.
- Southwell, C., J. Clarke, K. Reid, G. Watters and D. Ramm. 2004. Review of the CEMP standard methods and their delivery to the CEMP database. Document *WG-EMM-04/70*. CCAMLR, Hobart, Australia.

Underwood, A.J. 1991. Beyond BACI: experimental designs for detecting human environmental impacts on temporal variations in natural populations. *Aust. J. Mar. Freshwat. Res.*, 42: 569–587.

Whitehead, M.D. 1989. Sensitivity analysis for parameters of predatory species. CCAMLR Ecosystem Monitoring Program. In: *Selected Scientific Papers, 1989 (SC-CAMLR-SSP/6)*. CCAMLR, Hobart, Australia: 411–432.

Table 1: Variance components in the guard and crèche stages, with percentage of total variance in brackets, estimated from pre-impact data.

Stage	Component			
	Season	Pair	Penguin	Trip
Guard	0.038 (8)	0.039 (8)	0.118 (24)	0.289 (60)
Crèche	0.023 (2)	-	0.131 (14)	0.772 (83)

Table 2: Estimated variances and percentage reduction in variance (in brackets) as the number of penguins and number of trips per penguin are increased above a minimal combination of 10 penguins making one trip each. Results are based on calculation of variances as described in the text plus the variance component estimates in Table 1.

Stage	Number of penguins	Number of trips				
		1	2	3	4	8
Guard	10	0.087 (0)	0.072 (17)	0.067 (22)	0.065 (25)	0.061 (29)
	30	0.054 (37)	0.049 (43)	0.048 (45)	0.047 (46)	0.046 (47)
	50	0.048 (45)	0.045 (48)	0.044 (49)	0.043 (50)	0.043 (51)
	70	0.045 (48)	0.043 (50)	0.042 (51)	0.042 (52)	0.041 (52)
	100	0.043 (50)	0.041 (52)	0.041 (53)	0.041 (53)	0.040 (53)
Crèche	10	0.113 (0)	0.075 (34)	0.062 (45)	0.055 (51)	0.046 (60)
	30	0.053 (53)	0.040 (64)	0.036 (68)	0.034 (70)	0.031 (73)
	50	0.041 (64)	0.033 (71)	0.031 (73)	0.029 (74)	0.028 (76)
	70	0.036 (68)	0.030 (73)	0.029 (75)	0.028 (76)	0.026 (77)
	100	0.032 (72)	0.028 (75)	0.027 (76)	0.026 (77)	0.025 (78)

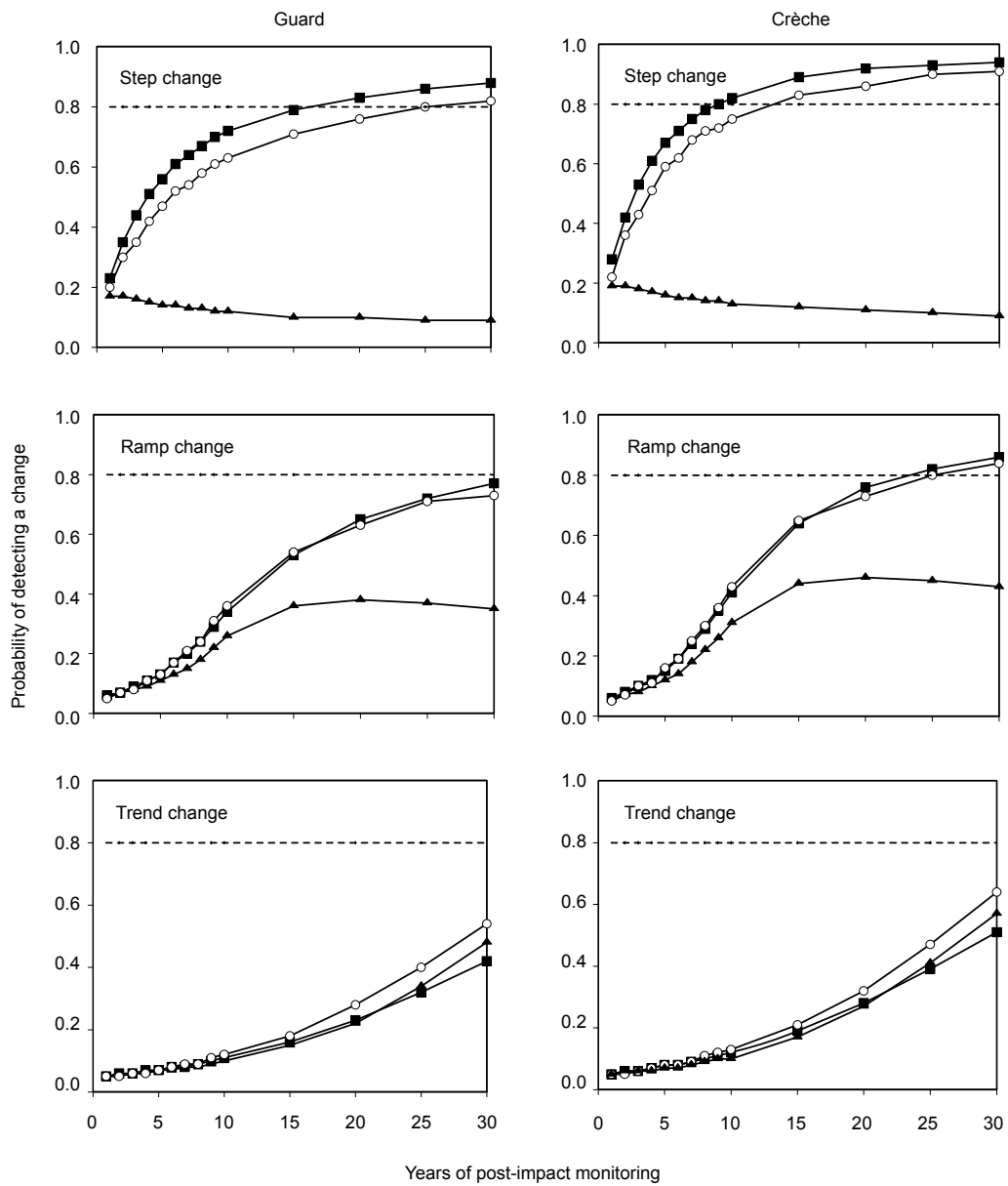


Figure 1: Probability of detecting a systematic increase in average foraging trip duration under a number of impact (step, ramp and trend increases of 25%) and monitoring (1–10, 15, 20, 25 and 30 years of post-impact monitoring; Type I error rate of 0.05) scenarios for the guard and crèche stages, given 12 years of pre-impact baseline data and using difference (■), slope (▲) and joint (○) statistics. The dashed horizontal line indicates power = 0.80.

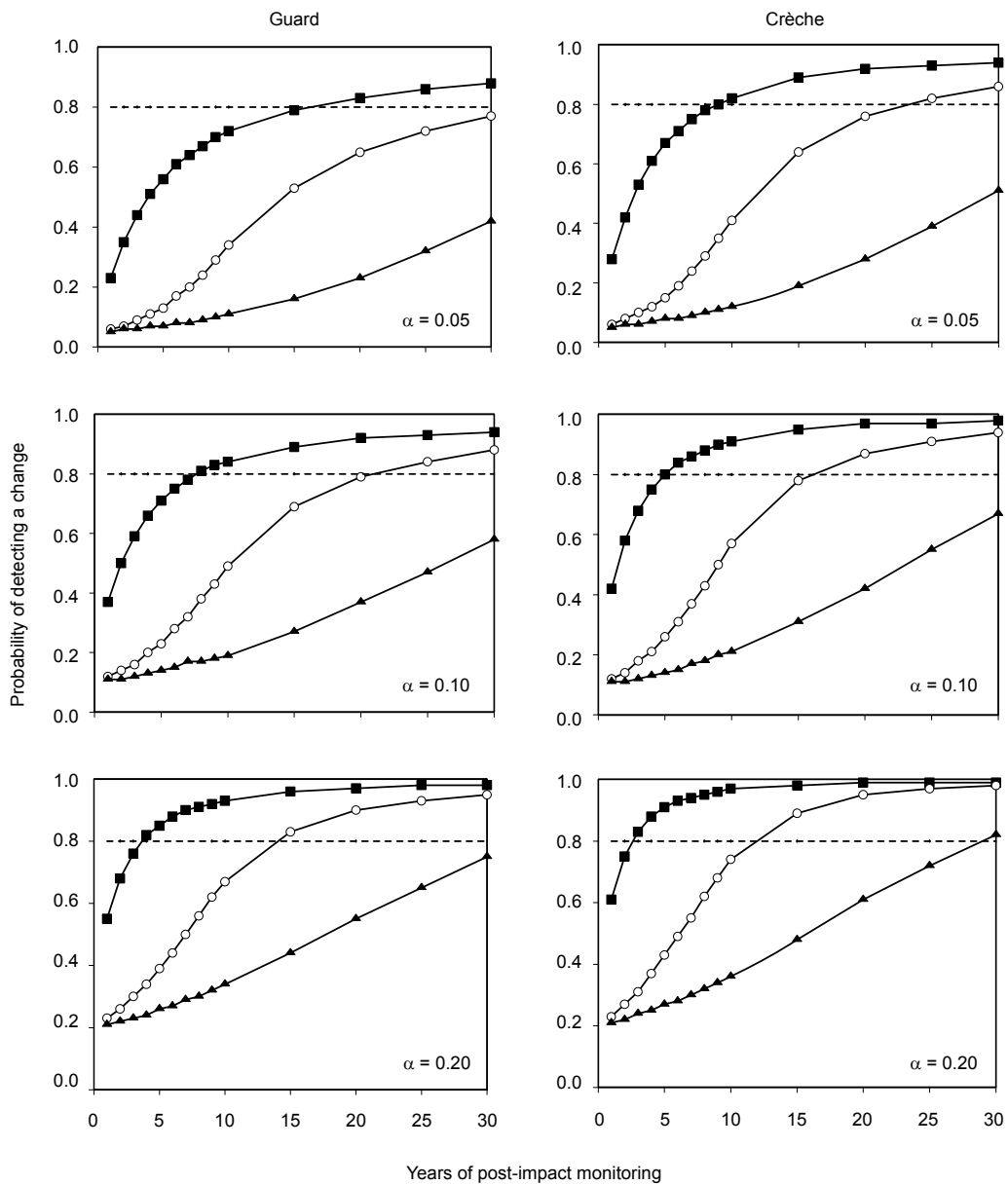


Figure 2: Probability of detecting a systematic increase in average foraging trip duration under a number of impact (step (■), ramp (○) and trend (▲) increases of 25%) and monitoring (1–10, 15, 20, 25 and 30 years of post-impact monitoring; Type I error rates of 0.05, 0.10 and 0.20) scenarios for the guard and crèche stages, given 12 years of pre-impact baseline data and using the difference statistic. The dashed horizontal line indicates power = 0.80.

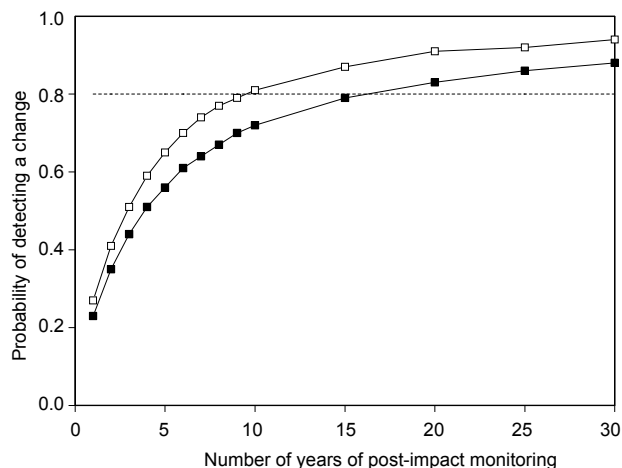


Figure 3: Probability of detecting a step increase of 25% in foraging trip duration during the guard stage without ice (■) and with ice (□) as a covariate, given 12 years of pre-impact baseline data. Tests are one-tailed with $\alpha = 0.05$, and power is calculated using the difference statistic. The dashed horizontal line indicates power = 0.80.

Liste des tableaux

- Tableau 1: Composantes de la variance aux stades de garde et de crèche avec, entre parenthèses, le pourcentage de la variance totale estimée à partir des données pré-impact.
- Tableau 2: Estimation des variances et, entre parenthèses, du pourcentage de réduction de la variance lorsque le nombre de manchots et le nombre de sorties par manchot sont augmentés au-delà d'un nombre minimal combiné de 10 manchots effectuant une sortie chacun. Les résultats sont fondés sur le calcul des variances décrit dans le texte et sur l'estimation des composantes de la variance figurant au tableau 1.

Liste des figures

- Figure 1: Probabilité de détection d'une augmentation systématique de la durée moyenne d'une sortie alimentaire dans divers scénarios d'impact (augmentation de 25% : marche, rampe et tendance) et de suivi (1-10, 15, 20, 25 et 30 ans de suivi post-impact; taux d'erreur de type I de 0,05) pour les stades de garde et de crèche, pour 12 années de données de base antérieures à l'impact et en utilisant les statistiques de différence (■), de pente (▲) et combinées (○). La ligne horizontale pointillée indique une puissance = 0,80.
- Figure 2: Probabilité de détection d'une augmentation systématique dans la durée moyenne d'une sortie alimentaire sous divers scénarios d'impact (augmentation de 25% sous forme de marche(■), de rampe (○) ou de tendance(▲)) et de suivi (1-10, 15, 20, 25 et 30 ans de suivi après l'impact; taux d'erreur de type I de 0,05, 0,10 et 0,20) pour les stades de garde et de crèche, pour 12 années de données de base antérieures à l'impact et en utilisant les statistiques de différence. La ligne horizontale pointillée indique une puissance = 0,80.
- Figure 3: Probabilité de détection d'une augmentation de 25%, sous forme de marche, dans la durée moyenne d'une sortie alimentaire pendant le stade de garde, sans glace (■) et avec glace (□) en tant que covariante, pour 12 années de données de base antérieures à l'impact. Les tests sont des tests à une queue, lorsque $\alpha = 0.05$ et la puissance est calculée au moyen des statistiques de différence. La ligne horizontale pointillée indique une puissance = 0,80.

Список таблиц

- Табл. 1: Компоненты дисперсии для ясельной стадии и периода присмотра; в скобках показан процент общей дисперсии, рассчитанный по данным до воздействия.

Табл. 2: Оценки дисперсии и процентное уменьшение дисперсии (в скобках) по мере увеличения числа пингвинов и числа походов на пингвина выше минимальной комбинации (10 пингвинов, каждый из которых совершает один поход). Результаты основаны на расчете дисперсии, как описано в тексте, и на оценках компонентов дисперсии в табл. 1.

Список рисунков

Рис. 1: Вероятность выявления систематического роста средней продолжительности похода за пищей в случае различных сценариев воздействия (скачкообразный, плавный и медленный рост на 25%) и мониторинга (1–10, 15, 20, 25 и 30 лет мониторинга после воздействия; величина ошибки первого рода 0.05) на ясельной стадии и в период присмотра; используются контрольные данные за 12 лет до воздействия и статистические показатели разницы (■), крутизны (▲) и их комбинации (d). Пунктирная горизонтальная линия показывает уровень = 0.80.

Рис. 2: Вероятность выявления систематического роста средней продолжительности похода за пищей в случае различных сценариев воздействия (скачкообразный (■), плавный (○) и медленный (▲) рост на 25%) и мониторинга (1–10, 15, 20, 25 и 30 лет мониторинга после воздействия; величина ошибки первого рода составляет 0.05, 0.10 и 0.20) на ясельной стадии и в период присмотра; используются контрольные данные за 12 лет до воздействия и статистический показатель разницы. Пунктирная горизонтальная линия показывает уровень = 0.80.

Рис. 3: Вероятность выявления скачкообразного изменения продолжительности похода за пищей на 25% в период присмотра, где ковариатой является отсутствие льда (■) и наличие льда (□); используются контрольные данные за 12 лет до воздействия. Критерии являются односторонними с $\alpha = 0.05$; вероятность рассчитывается по статистическому показателю разницы. Пунктирная горизонтальная линия показывает уровень = 0.80.

Lista de las tablas

Tabla 1: Componentes de la varianza en las etapas de cría y de guardería, indicándose entre paréntesis el porcentaje de la varianza total estimado a partir de los datos antes del impacto.

Tabla 2: Estimaciones de la varianza y porcentaje de reducción de la misma (entre paréntesis) a medida que aumenta el número de pingüinos y de los viajes de alimentación por pingüino, por sobre una combinación mínima de 10 pingüinos con un viaje de alimentación cada uno. Los resultados se basan en los cálculos de la varianza descritos en el texto más las estimaciones de los componentes de la varianza de la tabla 1.

Lista de las figuras

Figura 1: Probabilidad de detección de un aumento sistemático de la duración promedio de los viajes de alimentación bajo varias condiciones relativas al impacto (un aumento de 25% en el intervalo, la tendencia y la rampa) y al seguimiento (1–10, 15, 20, 25 y 30 años de seguimiento tras el impacto; suponiendo que la tasa de error Tipo I es igual a 0.05), tanto para la etapa de cría como de guardería, con datos recopilados durante 12 años antes del impacto para la línea de base y utilizando las pruebas estadísticas basadas en la diferencia (■), en la pendiente (▲) y en ambas (○). La línea entrecortada horizontal indica una potencia de = 0.80.

Figura 2: Probabilidad de detección de un aumento sistemático de la duración promedio de los viajes de alimentación bajo varias condiciones relativas al impacto (un aumento de 25% en el intervalo (■), la rampa (○) y la tendencia (▲) y al seguimiento (1–10, 15, 20, 25 y 30 años de seguimiento tras el impacto; suponiendo tasas del error Tipo I igual a 0.05, 0.10 y 0.20), tanto para la etapa de cría como de guardería, con datos recopilados durante 12 años antes del impacto para la línea de base y utilizando la prueba estadística basada en la diferencia. La línea entrecortada horizontal indica una potencia de = 0.80.

Figura 3: Probabilidad de detectar un aumento de 25% (en forma de intervalo) de la duración de los viajes de alimentación durante la etapa de cría, con (■) y sin (□) el factor hielo como covariante, con datos recopilados durante 12 años antes del impacto para la línea de base. Las pruebas son con un extremo (o unilaterales) con un valor de $\alpha = 0.05$, y la potencia se calcula mediante la prueba estadística basada en la diferencia. La línea entrecortada horizontal indica una potencia de = 0.80.

