# Accounting for Location Error in Kalman Filters: Integrating Animal Borne Sensor Data into Assimilation Schemes

**Aritra Sengupta[1]\*, Scott D. Foster[2], Toby A. Patterson[3], Mark Bravington[2]**

**1** Department of Statistics, The Ohio State University, Columbus, Ohio, United States of America, **2** CSIRO Mathematical and Information Sciences, Hobart, Tasmania, Australia, **3** CSIRO Wealth from Oceans Research Flagship, Castray Esplanade, Hobart, Tasmania, Australia

## Abstract

Data assimilation is a crucial aspect of modern oceanography. It allows the future forecasting and backward smoothing of ocean state from the noisy observations. Statistical methods are employed to perform these tasks and are often based on or related to the Kalman filter. Typically Kalman filters assumes that the locations associated with observations are known with certainty. This is reasonable for typical oceanographic measurement methods. Recently, however an alternative and abundant source of data comes from the deployment of ocean sensors on marine animals. This source of data has some attractive properties: unlike traditional oceanographic collection platforms, it is relatively cheap to collect, plentiful, has multiple scientific uses and users, and samples areas of the ocean that are often difficult of costly to sample. However, inherent uncertainty in the location of the observations is a barrier to full utilisation of animal-borne sensor data in data-assimilation schemes. In this article we examine this issue and suggest a simple approximation to explicitly incorporate the location uncertainty, while staying in the scope of Kalman-filter-like methods. The approximation stems from a Taylor-series approximation to elements of the updating equation.

**Competing Interests:** The authors have declared that no competing interests exist.

\* E-mail: sengupta.11@osu.edu

## Introduction

The process of updating physical ocean models using observations, to obtain accurate estimates of ocean state is referred to as data assimilation (DA) and is used to forecast current and future ocean conditions, as well as for hind-casting (backward smoothing) of historical states (e.g., [1], [2]). Data assimilation schemes must be computationally cheap, as the scale of oceanographic and atmospheric systems are generally large, typically with fine granularity in time, and large number of spatial cells. The ocean and atmosphere are continuously changing, so it is desirable to efficiently update model predictions (forecasts and hind-casts) with new data as it comes online. Also, the DA scheme needs to be able to use a wide variety of different data sources. Many examples of atmospheric and oceanographic models exist (for more details on data assimilation schemes see [1] and [3]). In this study we consider the use of spatially imprecise measurements in DA schemes – accurate measurements on the observed state variables, with imprecise spatial locations.

Modern biologging technology has brought a glut of observations of ocean temperature and salinity at depth (e.g., [4]), but a significant barrier to uptake of such data in physical models is the issue of spatial uncertainty. Traditionally, most data used in DA schemes are obtained from specially designed sampling devices, such as Argos floats (http://www.argo.ucsd.edu/), ship-based

instruments, Lagrangian drifters and remote sensing data. These platforms provide highly accurate information but are costly to deploy. The Lagrangian drifters have been used for tracking upper ocean water circulation and sea surface temperature (e.g., [5], [6]). However, these sampling platforms tend to under-sample in some areas; for example, Argo floats are often advected away from coastal areas, or are blocked from ice prone areas (e.g., [7], [8]). In contrast, a recent and cost effective addition to the suite of ocean sampling platforms has been miniaturized instruments attached to marine animals; the sensors sample depth/pressure, ambient temperature and conductivity (e.g. [9]). Sampling rates vary between instruments, but recent monitoring devices can collect data at 1 Hz, although operational sampling rates for lengthy deployments are often as low as 1/60 Hz.

Using free-swimming animals as data collection sources has many advantages. However, the key drawback with using this source of data is that the precise location of the observations are sometimes poorly known (e.g., [10]). Location data from oceanographic drifters and Argos floats, being derived from similar observation technology, also suffer from the problems of spatial error. However, for these technologies, there is enough spatial data associated with each ocean-state measurement so that the locations which are deemed inaccurate can be discarded, or straightforwardly corrected. With data collected from animal-borne sensors this is typically not the case, and its use in DA must

take this into account. This article demonstrates how incorporation of this sort of data into DA schemes is feasible, despite the associated location uncertainty.

Broadly speaking, there are two types of widely used animal-borne sensors, known as 'tags' : 1) data-storage (also called archival tags, see [11], [12]) for which non-spatial sensor data is used to spatially locate the animal (see [10], [13–14]) and, 2) satellite tags which are spatially located by satellite providers such as CLS-Argos (see www.cls.fr).
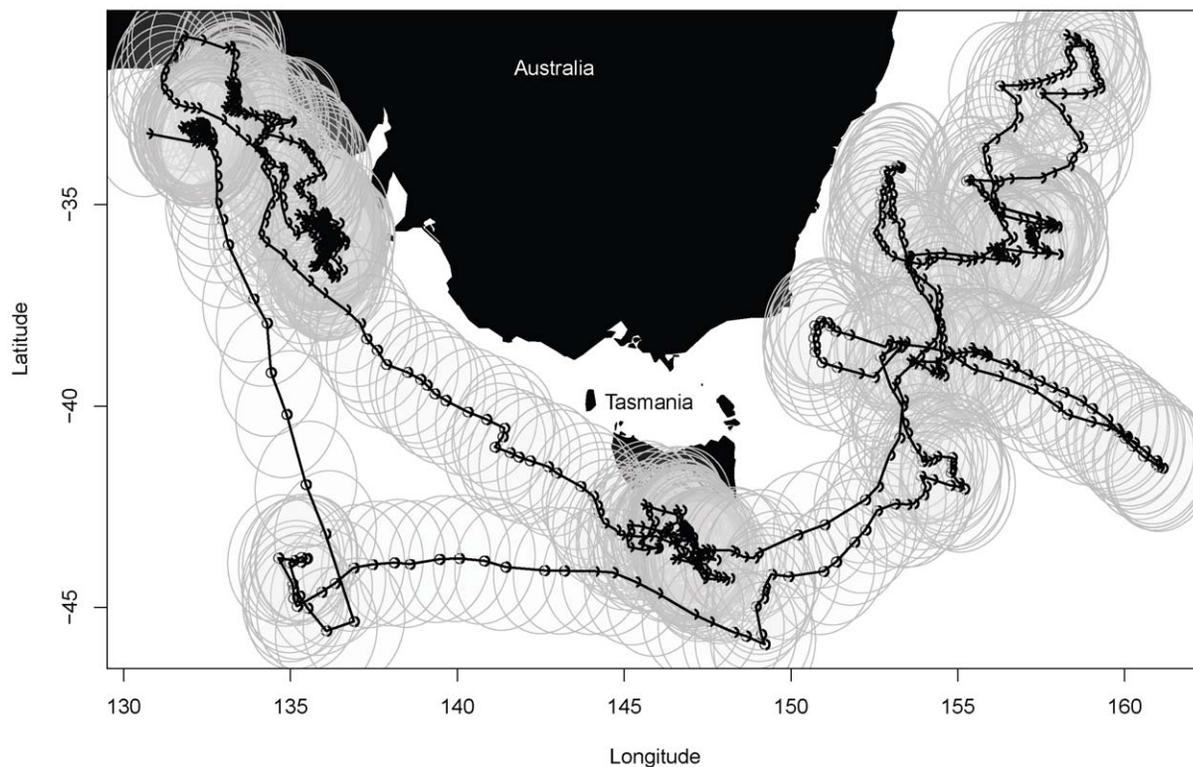
Satellite tags are attached externally to animals which spend sufficient time at the surface for the tag antenna to be exposed, and thus able to transmit. Therefore, satellite tags send near-real time data. Satellite tags are often more expensive, but have the advantage of far lower degrees of spatial error in positions (e.g., [15] and [16]). However, most satellite tags need to drastically summarize raw sensor data streams because of the low bandwidth and limited battery life available for data transmission (e.g., [17]). In the context of a DA scheme, this means that data from satellite tags could be used for forward- and hind-casting.

In contrast, data-storage tags, typically used for tracking non-air breathing animals, simply store data on board, and data retrieval relies on the tags being retrieved when animals are recaptured. Such data are highly detailed with many thousands of observations, but are geo-located only infrequently (e.g., only twice per 24 h) and with low spatial accuracy (see example in Figure 1). In these data, latitudinal errors are typically much greater than longitudinal errors, and vary systematically through time (e.g., [18]). In the DA scheme context, these data sources are therefore retrospective, and thus are primarily useful for hind-casting.

Data assimilation for many oceanographic models is based on variants of the Kalman filter algorithm (see [19] for a statistical context), with hind-casting naturally performed by corresponding variants of Kalman smoothing (e.g., [20]). In this article we describe some simple adjustments to the Kalman filter algorithm which allow the use of spatially imprecise data in the DA scheme. We then briefly consider the applicability of these adjustments to more complex DA schemes. Note that in this article, a direct and efficient approach that could be implemented within the extensive existing libraries of oceanographic modelling code, is the key, and thus we avoid consideration of Bayesian hierarchical models which tend to be computationally prohibitive in operational DA schemes.

We proceed as follows: A brief review of the Kalman filter models and the updating equations is given. Then we detail adjustments to the Kalman filter procedure which account for location errors. We also discuss the Kalman smoothing algorithm in the presence of location errors. The setup and the results of some simple simulation studies are discussed. Finally, we discuss our methods, the scope and limitations of this study, and some of the possible extension of the results. Throughout this article we will assume that the location error variance is known for all time points although, in practice these need to be estimated. However, they may be directly obtained from dedicated studies (e.g., [21]).

In this short article, our goal is not to present a fully operational oceanographic DA scheme that handles location error; that would be a much larger task. Our contribution is to demonstrate, via application to simulated data, that accounting for spatial uncertainty is a surmountable problem, without wholesale adjustment to standard techniques.



**Figure 1. An example of movement data obtained from a data storage tag deployed on a tuna (CSIRO unpublished data).** The Location estimates were derived from the state-space model approach developed by [25]. This method yields a point-estimate along with error-variances from the 95% error ellipses depicted here were derived. Associated with this track are records of temperature-at-depth recorded every minute over a period of approximately 12 months (data not shown).
doi:10.1371/journal.pone.0042093.g001

## Materials and Methods

### The Kalman Filter

A brief review of the Kalman filter is now given. For a more thorough treatment see [19] or [22]. We use the notation used in [19] but acknowledge that this is not the only choice. A notable alternative is [23]. In Appendix S1 we give a glossary of terms and a description of each one, for both sets of terminologies. This should aid translation for those familiar with [23].

Let $\mathbf{Y}_1, \ldots \mathbf{Y}_t$ denote the observed values of a variable of interest at time points $1, 2, \ldots, t$ respectively. The observations may be vector or scalar depending on the particular system under study in the DA setup, and the spatial locations of these observations are given by $\mathbf{X}_1, \ldots, \mathbf{X}_t$. We assume that $\mathbf{Y}_t$ depends on an unobserved quantity $\mathbf{Z}_t$, known as the state of the nature, or the system state. Typically, the value of $\mathbf{Y}_t$ also depends on the location of the observation, $\mathbf{X}_t$. The model for $[\mathbf{Y}_t | (\mathbf{X}_{\leq t}, \mathbf{Z}_{\leq t})]$, which is called the *observation equation*, is

$$\mathbf{Y}_t | (\mathbf{X}_{\leq t}, \mathbf{Z}_{\leq t}) = \mathbf{F}_t(\mathbf{X}_t)\mathbf{Z}_t + \mathbf{d}_t + \mathbf{v}_t, \tag{1}$$

where $\mathbf{F}_t(\mathbf{X}_t)$ is a known quantity which may change with time and measurement location. The model for $\mathbf{Y}_t$ may be governed by its correlation with the components of $\mathbf{Z}_t$. In such cases, the elements of $\mathbf{F}_t(\cdot)$ will depend on the covariance between $\mathbf{Y}_t$ and $\mathbf{Z}_t$. In (1), $\mathbf{d}_t$ is a known vector of the same dimension as $\mathbf{Y}_t$ and it may, or may not change with time. The *observation error*, $\mathbf{v}_t$, is assumed to be normally distributed with mean zero and a known variance $\mathbf{V}_t$, which may also be time dependent, i.e., $\mathbf{v}_t \sim N(0, \mathbf{V}_t)$.

The model for the *state of nature*, $\mathbf{Z}_t$, is given by the *system equation* which is of the form

$$\mathbf{Z}_t = \mathbf{G}_t \mathbf{Z}_{t-1} + \mathbf{c}_t + \mathbf{w}_t, \tag{2}$$

where $\mathbf{G}_t$ is a known quantity, $\mathbf{c}_t$ is some known vector of the same dimension as $\mathbf{Z}_t$, and the *system equation error*, $\mathbf{w}_t \sim N(0, \mathbf{W}_t)$, where $\mathbf{W}_t$ is assumed known. Here also $\mathbf{G}_t$, $\mathbf{c}_t$ and $\mathbf{W}_t$ may or may not change with time. The system equation (2) does not depend on the location of $\mathbf{Y}_t$. If the $\mathbf{Z}_t$ vectors are measured over space, their locations are pre-determined, and not subject to any measurement error. We also assume that $\mathbf{w}_t$ and $\mathbf{v}_t$ are independent.

One can write down the joint distribution $[\mathbf{Z}_t, \mathbf{Y}_t]$, conditional on $\mathbf{Y}_{\leq t-1}$ and $\mathbf{X}_{\leq t}$ as:

$$\left( \begin{pmatrix} \mathbf{Z}_t \\ \mathbf{Y}_t \end{pmatrix} | \mathbf{X}_{\leq t}, \mathbf{Y}_{\leq t-1} \right) \sim N \left( \begin{pmatrix} \hat{\mathbf{Z}}_{t|t-1} \equiv \mathbf{G}_t \hat{\mathbf{Z}}_{t-1|t-1} + \mathbf{c}_t \\ \hat{\mathbf{Y}}_t \end{pmatrix}, \right.$$
$$\left. \begin{pmatrix} \mathbf{R}_t & \mathbf{R}_t \mathbf{F}_t(\mathbf{X}_t)^{\mathrm{T}} \\ \mathbf{F}_t(\mathbf{X}_t)\mathbf{R}_t & \mathbf{V}_t + \mathbf{F}_t(\mathbf{X}_t)\mathbf{R}_t\mathbf{F}_t(\mathbf{X}_t)^{\mathrm{T}} \end{pmatrix} \right) \tag{3}$$

where $\hat{\mathbf{Z}}_{t-1|t-1}$ is the estimate of the system state at time point $(t-1)$ using all available information upto time $(t-1)$, $\hat{\mathbf{Y}}_t \equiv \mathbf{F}_t(\mathbf{X}_t)(\mathbf{G}_t\hat{\mathbf{Z}}_{t-1|t-1} + \mathbf{c}_t) + \mathbf{d}_t$, and $\mathbf{S}_{t|t-1} \equiv \mathbf{R}_t \equiv (\mathbf{G}_t\mathbf{S}_{t-1|t-1}\mathbf{G}_t^{\mathrm{T}} + \mathbf{W}_t)$ (see [19]).

Once we obtain an observation on $\mathbf{Y}_t$, we can use the joint distribution (3) to obtain the expectation and variance for the posterior distribution $[\mathbf{Z}_t | \mathbf{Y}_{\leq t}, \mathbf{X}_{\leq t}]$. The expected value and the variance of the posterior distribution are:

$$\hat{\mathbf{Z}}_{t|t} \equiv (\mathbf{G}_t\hat{\mathbf{Z}}_{t-1|t-1} + \mathbf{c}_t) + \mathbf{R}_t\mathbf{F}_t(\mathbf{X}_t)^{\mathrm{T}} \left( \mathbf{V}_t + \mathbf{F}_t(\mathbf{X}_t)\mathbf{R}_t\mathbf{F}_t(\mathbf{X}_t))^{\mathrm{T}} \right)^{-1}$$
$$\left( \mathbf{Y}_t - \mathbf{F}_t(\mathbf{X}_t)\left(\mathbf{G}_t\hat{\mathbf{Z}}_{t-1|t-1} + \mathbf{c}_t\right) - \mathbf{d}_t \right) \tag{4}$$

and,

$$\mathbf{\Sigma}_{t|t} \equiv \mathbf{R}_t - \mathbf{R}_t\mathbf{F}_t(\mathbf{X}_t)^{\mathrm{T}} \left( \mathbf{V}_t + \mathbf{F}_t(\mathbf{X}_t)\mathbf{R}_t\mathbf{F}_t(\mathbf{X}_t)^{\mathrm{T}} \right)^{-1} \mathbf{F}_t(\mathbf{X}_t)\mathbf{R}_t. \tag{5}$$

Note that the posterior distribution depends on the location $\mathbf{X}_t$.

### Adjustments for Location Error

Suppose we have an estimate for the location of an observation $\mathbf{Y}_t$, available from auxiliary data (i.e. estimates of location derived from the tag data). Let us denote the estimate by $\xi_t$. We assume that $\mathbf{X}_t$, the true location of the observation, is a random variable, the distribution of which is centered around the estimate $\xi_t$ and has some variance which will be denoted by $V(\mathbf{X}_t | \xi_t)$. That is, we have

$$\mathbf{X}_t \sim N(\mathbf{x}_t, V(\mathbf{X}_t | \mathbf{x}_t)). \tag{6}$$

We will assume that we have an estimate $V(\mathbf{X}_t | \xi_t)$ through previous experimental data (e.g. [15]). We further assume that $\mathbf{X}_t$ is independent between time steps. The necessary adjustments to the filtering algorithm will now be made using a first-order Taylor-series expansion.

#### Adjustments to the Joint Distribution. $[\mathbf{Z}_t, \mathbf{Y}_t | \mathbf{Y}_{\leq t-1}]$

Using a first-order approximation, the expectation of $\mathbf{Y}_t$ is (see Appendix S2 for all the derivations in this section)

$$\hat{\mathbf{Y}}_t \equiv E(\mathbf{Y}_t | \xi_t, \mathbf{Y}_{\leq t-1}) \approx E(\mathbf{Y}_t | \mathbf{X}_t = \xi_t, \mathbf{Y}_{\leq t-1}). \tag{7}$$
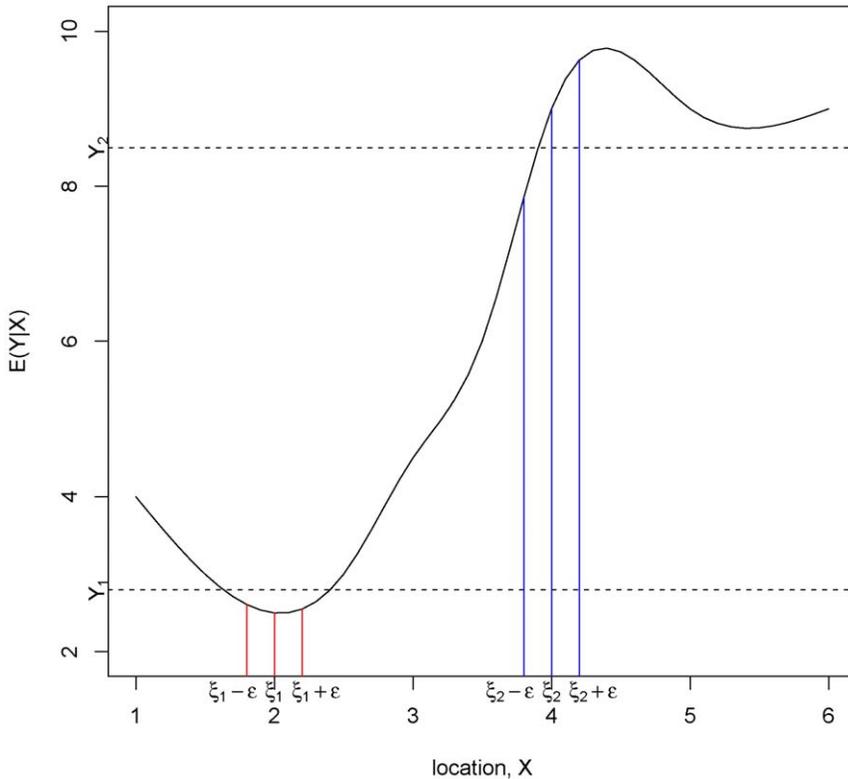
where the conditioning on the left hand side is with respect to (w.r.t.) the noisy location, and the conditioning on the right hand side is done w.r.t. the true location, treating the noisy estimate as the truth.

This suggests that using the noisy estimate of location is a valid approximation. However, there is potentially substantial bias in the variance as

$$V(\mathbf{Y}_t | \xi_t, \mathbf{Y}_{\leq t-1}) \approx V(\mathbf{Y}_t | \mathbf{X}_t = \xi_t, \mathbf{Y}_{\leq t-1})$$
$$+ \left\{ \left[ \frac{\partial}{\partial \mathbf{X}_t}(E(\mathbf{Y}_t | \mathbf{X}_t, \mathbf{Y}_{\leq t-1})) \right]_{\mathbf{X}_t = \xi_t}^{\mathrm{T}} \right.$$
$$V(\mathbf{X}_t | \xi_t) \left[ \frac{\partial}{\partial \mathbf{X}_t}(E(\mathbf{Y}_t | \mathbf{X}_t, \mathbf{Y}_{\leq t-1})) \right]_{\mathbf{X}_t = \xi_t}. \tag{8}$$

Reasuringly, the approximate variance in the unknown location situation is inflated as compared to knowing locations.

The term $\frac{\partial}{\partial \mathbf{X}_t}(E(\mathbf{Y}_t | \mathbf{X}_t, \mathbf{Y}_{\leq t-1}))$ measures how much the expectation will change for a small change in the location $\mathbf{X}_t$. The variance inflation term in (8) will not be too significant if the local slope of the surface of the expectation around $\mathbf{X}_t = \xi_t$ is small (see Figure 2). In such a situation, the precise location of the observation is less influential. However, if the local slope is large, then the precise location does matter. See Figure 2 for a pictorial illustration.

**Figure 2. A diagrammatic representation of the problem where the curve corresponds to** $E(\mathbf{Y}_t|\mathbf{X}_t,\mathbf{Y}_{\leq t-1})$. $\xi_1$ is a point with small local slope (in this case, the precision of the location is less influential) and $\xi_2$ is a point with a larger local slope (here, the precise location is influential).
doi:10.1371/journal.pone.0042093.g002

We complete the adjustments to the joint distribution by considering the covariance between $\mathbf{Y}_t$ and $\mathbf{Z}_t$,

$$\mathrm{Cov}(\mathbf{Y}_t,\mathbf{Z}_t|\mathbf{Y}_{\leq t-1},\xi_t) \approx Cov(\mathbf{Y}_t,\mathbf{Z}_t|\mathbf{Y}_{\leq t-1},\mathbf{X}_t=\xi_t)$$
$$= \mathbf{R}_t \mathbf{F}_t(\xi_t). \tag{9}$$

So, the covariance, like the expectation, requires no adjustment.

Combining the expectation, the variance, and the covariance, the joint distribution of $\mathbf{Z}_t$ and $\mathbf{Y}_t$, conditional on $\mathbf{Y}_{\leq t-1}$ and $\xi_t$ can be written as,

$$\left( \binom{\mathbf{Z}_t}{\mathbf{Y}_t} \bigg| \xi_t, \mathbf{Y}_{\leq t-1} \right) \sim N\left( \binom{\hat{\mathbf{Z}}_{t|t-1} \equiv \mathbf{G}_t\hat{\mathbf{Z}}_{t-1|t-1}+\mathbf{c}_t}{\hat{\mathbf{Y}}_t}, \right.$$
$$\left. \begin{pmatrix} \mathbf{R}_t & \mathbf{R}_t\mathbf{F}_t(\xi_t)^{\mathrm{T}} \\ \mathbf{F}_t(\xi_t)\mathbf{R}_t & V(\mathbf{Y}_t|\xi_t,\mathbf{Y}_{\leq t-1}) \end{pmatrix} \right), \tag{10}$$

where $V(\mathbf{Y}_t|\xi_t,\mathbf{Y}_{\leq t-1})$ is given by (8) and $\hat{\mathbf{Y}}_t$ is given by (7). This should be contrasted with (3). The variance of $\mathbf{Z}_t$ is larger than that in (3).

The Taylor-series expansion shown above retained terms up to the first-order. Better approximations may be obtained considering the higher order terms. In Appendix S3 we derive the second-order correction for the expectation. However, computation of these extra terms is more expensive and difficult to code, so we do not pursue it further in this article.

**The Posterior Estimates.** Recall that the posterior estimate for the state of the system at time point $(t-1)$ is $\hat{\mathbf{Z}}_{t-1|t-1}$, which is obtained using all the information available up to time point

$(t-1)$. Then we can write down the joint distribution of $\mathbf{Z}_t$ and $\mathbf{Y}_t$, conditional on $\mathbf{Y}_{\leq t-1}$ and $\xi_t$, as in (10). From this joint normal distribution, we can derive the posterior distribution $[\mathbf{Z}_t|(\mathbf{Y}_{\leq t},\xi_t)]$ (see [20]).

The mean of the posterior distribution is

$$\hat{\mathbf{Z}}_{t|t} \equiv \left( \mathbf{G}_t\hat{\mathbf{Z}}_{t-1|t-1}+\mathbf{c}_t \right) + \mathbf{R}_t\mathbf{F}_t(\xi_t)^{\mathrm{T}}\left( (V(\mathbf{Y}_t|\xi_t,\mathbf{Y}_{\leq t-1}))^{-1} \right.$$
$$\left. \{\mathbf{Y}_t - \mathbf{F}_t(\xi_t)(\mathbf{G}_t\hat{\mathbf{Z}}_{t-1|t-1}+\mathbf{c}_t)-\mathbf{d}_t\}, \right. \tag{11}$$

and the variance is

$$\mathbf{\Sigma}_{t|t}^{*} \equiv \mathbf{R}_t - \mathbf{R}_t\mathbf{F}_t(\xi_t)(V(\mathbf{Y}_t|\xi_t,\mathbf{Y}_{\leq t-1}))^{-1}\mathbf{F}_t(\xi_t)\mathbf{R}_t, \tag{12}$$

where $V(\mathbf{Y}_t|\xi_t,\mathbf{Y}_{\leq t-1})$ is given by (8). The posterior expectation and variance are both affected by location error through the alteration of $V(\mathbf{Y}_t|\xi_t,\mathbf{Y}_{\leq t-1})$. This posterior distribution gives our state of knowledge about $\mathbf{Z}_t$ at time $t$, using all the information available up to that time.

## Kalman Smoothing

Until now, we have considered only Kalman filtering, where an estimate of a signal $\mathbf{Z}_t$ was made from considering all the previous observations. The process will produce the *best forecasts* but it will not produce the *best hindcast* (estimate of the entire time series). To obtain the best hindcasts we need to consider all the data, both previous and future. We denote this estimate as $\hat{\mathbf{Z}}_{t|T} \equiv E(\mathbf{Z}_t|\mathbf{Y}_1,\dots,\mathbf{Y}_T)$ and refer to it as the Kalman smoothed estimate. We note that for many tag types (e.g. archival tags) the

data from animal tags can only be used for hind-casting as the data cannot be made available in time for forecasting. This is due to the problems of tag and data retrieval, and data manipulation. It is possible that this process can be sped up in future.

A detailed discussion and the derivation of the Kalman smoothing equations, in the general setting, can be found in [20]. In our current work we are concerned with the so-called fixed-interval smoothing because of its relevance to hind-casting problems. Accounting for location error in the Kalman smoothing algorithm can be obtained directly by following the steps in [20], using the equations in the previous section.

The approach is to run a Kalman filter forward in time over the full interval $[0,T]$, storing state estimates at each update, and then these stored quantities are run backwards in time to obtain the smoothed estimates. This process generates the smoothed estimates in reverse sequence, $\hat{\mathbf{Z}}_{T-1|T}, \hat{\mathbf{Z}}_{T-2|T}, \ldots \hat{\mathbf{Z}}_{1|T}$. The adjusted Kalman smoothing equations are

$$
\begin{aligned}
\hat{\mathbf{Z}}_{j-1|T} &= \hat{\mathbf{Z}}_{j-1|j-1} + \mathbf{S}_{j-1|j-1} \mathbf{G}_j^{\mathrm{T}} \mathbf{\Sigma}_{j|j-1}^{-1} \left( \hat{\mathbf{Z}}_{j|T} - \hat{\mathbf{Z}}_{j|j-1} \right) \\
&= \hat{\mathbf{Z}}_{j-1|j-1} + \mathbf{A}_{j-1} \left( \hat{\mathbf{Z}}_{j|T} - \hat{\mathbf{Z}}_{j|j-1} \right),
\end{aligned}
\tag{13}
$$

where $\mathbf{A}_{j-1} \equiv \mathbf{\Sigma}_{j-1|j-1} \mathbf{G}_j^{\mathrm{T}} \mathbf{\Sigma}_{j|j-1}^{-1}$. The recursion for the error covariance is

$$
\mathbf{\Sigma}_{j-1|T} = \mathbf{\Sigma}_{j-1|j-1} + \mathbf{A}_{j-1} \left( \mathbf{\Sigma}_{j|T} - \mathbf{\Sigma}_{j|j-1} \right),
\tag{14}
$$

where recall that $\hat{\mathbf{Z}}_{j|j-1} = \mathbf{G}_j \hat{\mathbf{Z}}_{j-1|j-1} + \mathbf{c}_t$, $\hat{\mathbf{S}}_{j|j-1} = \mathbf{R}_j$, and the estimates $\hat{\mathbf{Z}}_{j-1|j-1}$ and $\mathbf{\Sigma}_{j-1|j-1}$ are obtained using (11) and (12).

Thus, to account for location error in Kalman smoothing one needs to only adjust the forward moving filtering process. The fixed-interval smoothing equations have no further dependence on the location of the observation. The information in the observations has already been fully incorporated during the filtering process.

## Results

### One-Dimensional Simulation System

To test the efficacy of the adjustments for location error in Kalman filter and Kalman smoothing we performed a simple one dimensional simulation along the surface of a ring. In this simple simulation, we assume that our simulated animal moves along the surface of the ring taking measurements $\mathbf{Y}_t$ at time $t$ from 11 possible observation locations labelled from 0 to 10 (the setup of the simulation is shown in Figure 3). Let us pretend, for the sake of illustration, that this measurement is of temperature. We introduce a temperature gradient via *a source* and *a sink* at two particular locations as shown in Figure 3. At the source, for each time point, there is an increase of $1^\circ$C in temperature and the sink absorbs $1^\circ$C. The state vector of interest is therefore $\mathbf{Z} = (Z_0, Z_1, \ldots, Z_{10})^{\mathrm{T}}$, the true water temperature at each of the locations.

A continuous temperature gradient exists along the surface of the ring. However, the measured temperature of water at any given location was considered to be a linear interpolation of the two neighboring values of $\mathbf{Z}$, plus some random noise. The model for $\mathbf{Y}_t$ at any given time $t$, given the location of the observation $\mathbf{X}_t$, is given by

$$
\mathbf{Y}_t | \mathbf{X}_t, \mathbf{Z}_t = \omega_{1,t} Z_{\lfloor \mathbf{X}_t \rfloor, t} + \omega_{2,t} Z_{\lfloor \mathbf{X}_t \rfloor + 1, t} + \mathbf{v}_t,
\tag{15}
$$

where $\lfloor \mathbf{X}_t \rfloor$ denotes the largest integer less than $\mathbf{X}_t$ and the weights $\omega_1$ and $\omega_2$ are determined as $\omega_1 = \lfloor \mathbf{X}_t \rfloor + 1 - \mathbf{X}_t$ and $\omega_2 = \mathbf{X}_t - \lfloor \mathbf{X}_t \rfloor$. The expectation is a linear function of the true location $\mathbf{X}_t$. Note that $\omega_1$ and $\omega_2$ are not differentiable at the boundaries ($\mathbf{X}_t = 0, \ldots 10$). However,

$$
P(\mathbf{X}_t = n) = 0
$$

where, $n = 0, \ldots, 10$. Therefore, this should not be a problem. The observation error $\mathbf{v}_t$ was taken to be a $N(0, 0.01)$ variable.

The model for the evolution of $\mathbf{Z}$ over time was taken to be

$$
\mathbf{Z}_t = \mathbf{G} \mathbf{Z}_{t-1} + \mathbf{d} + \mathbf{w}_t
\tag{16}
$$

with

$$
\mathbf{G} = \begin{pmatrix}
0.5 & 0.25 & 0 & \ldots & 0.25 \\
0.25 & 0.5 & 0.25 & \ldots & 0 \\
\ldots & \ldots & \ldots & \ldots & \ldots \\
0.25 & 0 & \ldots & 0.25 & 0.5
\end{pmatrix}.
$$

and therefore represents a diffusion process. The constant $\mathbf{d}$ is the vector given by:

$$
\mathbf{d} = (0, 1, 0, 0, 0, 0, -1, 0, 0, 0, 0)^{\mathrm{T}}
$$

encapsulating the source and sink concept described earlier. The system equation error is $\mathbf{w}_t \sim N(0, \mathbf{W}_t)$, where $\mathbf{W}_t = \mathrm{diag}(0.1, 0.1, \ldots, 0.1)$.
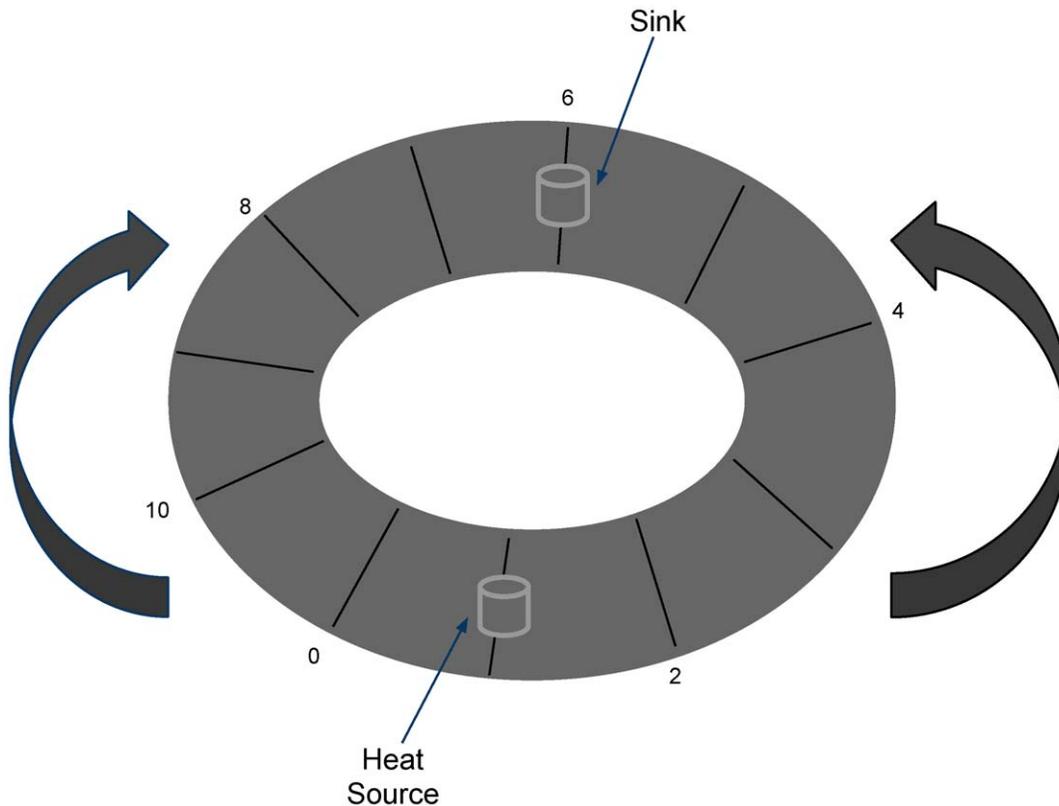
The filtering process for each data set was initialized with all state variables equal to $10^\circ$C. The system was evolved over time according to equation (16). The animal's location $\mathbf{X}_t$ was initialized to be at location 5 and evolves according a random walk, $\mathbf{X}_{t+1} \sim N(\mathbf{X}_t, 1)$. The $\mathbf{Y}_t$ vectors were simulated according to the model given by (15), for 100 time steps. A total of 1000 data sets were created using this model.

We compared several different variants of the Kalman filtering DA scheme. We estimated the state variables using the true locations and the noisy locations. Additionally, when using the noisy locations we used the filtering process ignoring the uncertainty in locations (eqn. (3)–(5)) and also the filtering process accounting for location uncertainty (eqn. (10)–(12)). Kalman smoothing was used to hind-cast the system states. The simulated values and the predictions from the filtering and smoothing, along with the mean square prediction errors (MSPEs), were recorded from 1000 simulation runs for each of the DA schemes. This was repeated for three different values of the location-error variance, namely 0.01, 0.1, and 1.

### One Dimensional Simulation Results

A summary of MSPEs for the different methods are given in Table 1. When the measurement error variance was small (0.01), there was little difference in the results for the three procedures. This is not surprising as the locations have low uncertainty. Increasing the location error variance decreases the performance of the standard Kalman filter. In these cases the location error adjusted Kalman filter updates performs substantially better (see Table 1).

Surprisingly, the Kalman smoothing algorithm applied to the data with noisy locations sometimes achieved a worse MSPE than the comparable Kalman filtering algorithm, when we didn't account for the noise while moving forward in time (i.e., when we

**Figure 3. Figure showing the simulation set-up.** Heat flows in the direction from the *Source* to the *Sink*.
doi:10.1371/journal.pone.0042093.g003

used the standard Kalman filter algorithm in the presence of location errors). This is contrary to prior expectation from theoretical considerations. It implies that the underlying statistical model used in Kalman filtering does not represent the data adequately in this situation. On the other hand, applying the smoothing algorithm to the estimates obtained using the adjusted Kalman filter algorithm led to a slight improvement in the results after smoothing. This is expected and indicates that the underlying statistical model is more consistent than that which ignores the location errors.

## Two-Dimensional Simulation on the Surface of a Torus

The one-dimensional simulation was extended to two dimensions by assuming a model spatial domain where both the X-coordinates and the Y-coordinates are joined at the ends to form a torus like structure. The X-axis coordinates ranged from $0-10$ and those along the Y-axis ranged from $0-12$. Simulations were run for 200 time steps and a heat source was located along the Y-coordinate 0 while a heat sink was located along the Y-coordinate 5.

For the two dimensional case, the observation equation considered was of the form

**Table 1.** Mean MSPE ($\pm$ standard deviation) of a 1000 simulated data sets with predictions from Kalman filtering and Kalman smoothing.

| | Location Error? | Measurement error variance | | |
|---|---|---|---|---|
| | | 0.01 | 0.1 | 1 |
| Usual KF | No | 0.237 ($\pm$0.014) | 0.238 ($\pm$0.014) | 0.237 ($\pm$0.014) |
| Smoothing | No | 0.185 ($\pm$0.009) | 0.186 ($\pm$0.009) | 0.185 ($\pm$0.009) |
| Usual KF | Yes | 0.255 ($\pm$0.016) | 0.442 ($\pm$0.032) | 3.515 ($\pm$0.320) |
| Smoothing | Yes | 0.202 ($\pm$0.010) | 0.392 ($\pm$0.032) | 3.892 ($\pm$0.396) |
| Adjusted KF | Yes | 0.251 ($\pm$0.016) | 0.313 ($\pm$0.023) | 0.538 ($\pm$0.078) |
| Smoothing | Yes | 0.198 ($\pm$0.010) | 0.253 ($\pm$0.017) | 0.423 ($\pm$ 0.061) |

Each data set contains 100 time steps and imitates an animal's movement over a one-dimensional ring. Location error is included and excluded (first rows) to gauge its effect.
doi:10.1371/journal.pone.0042093.t001

**Table 2.** Mean MSPE ($\pm$ standard deviation) of a 1000 simulated dat sets with predictions from Kalman filtering and Kalman smoothing.

| | MSPE ($\pm$ S.D) |
|---|---|
| Usual KF updates; no location error | 1.98 ($\pm$ 0.12) |
| Smoothing | 1.85 ($\pm$ 0.08) |
| Usual KF updates in presence of location errors | 3.86 ($\pm$ 0.71) |
| Smoothing | 4.28 ($\pm$ 0.94) |
| Adjusted KF updates | 2.27 ($\pm$ 0.25) |
| Smoothing | 2.18 ($\pm$ 0.21) |

Each data set contains 200 time steps and imitates an animal's movement over a torus. Location error is included and excluded (first rows) to gauge its effect.
doi:10.1371/journal.pone.0042093.t002

$$\mathbf{Y}_t | (\mathbf{X}_t = (x,y), \mathbf{Z}_t) = \omega_1 Z_{(\lfloor x \rfloor, \lfloor y \rfloor);t} + \omega_2 Z_{(\lfloor x \rfloor + 1, \lfloor y \rfloor);t} \\ + \omega_3 Z_{(\lfloor x \rfloor, \lfloor y \rfloor + 1);t} + \omega_4 Z_{(\lfloor x \rfloor + 1, \lfloor y \rfloor + 1);t} + \mathbf{v}_t \qquad (17)$$

where $\mathbf{X}_t = (x,y)$ denotes the true location of the measurement. The weights $\{\omega_1, \omega_2, \omega_3, \omega_4\}$ were determined according to bivariate linear interpolation (see [24]). The observation error $\mathbf{v}_t$ was taken to be a N(0,0.1) variable. The system equation to describe the evolution of the state over time was defined as

$$Z_{(x,y);t} = 0.4 Z_{(x,y);t-1} + 0.15 Z_{(x-1,y);t-1} + 0.15 Z_{(x+1,y);t-1} \\ + 0.15 Z_{(x,y-1);t-1} + 0.15 Z_{(x,y+1);t-1} + I(y=0) - I(y=5) + \mathbf{w}_t, \qquad (18)$$

where $I(\cdot)$ denotes the indicator function. The two terms involving the indicator function are the heat source and the sink. In these simulations the system equation error $\mathbf{w}_t$ was taken to be N(0,1). The location error variance was held at 1. A total of 1000 data sets were generated using this model and MSPE was used as the model's performance measure.

For the two dimensional simulation we again estimated the state vector using knowledge of both the true locations and the noisy locations. When using the noisy locations, we used both the usual filtering algorithm and the location error adjusted algorithm. As in the one-dimensional study, we applied both the filtering and the smoothing phase, again repeating the simulation/estimation procedure for 1000 iterations, setting the location error variance to one. A summary of the results are shown in Table 2.

In essence, the two-dimensional simulation results were concordant with the one dimensional simulation results. The adjusted Kalman filter updating equations again gave results that were better than those from an unadjusted DA scheme when there was uncertainty in location. In the presence of location error, the smoothing algorithm applied to the estimates obtained by applying the usual filtering process, produce estimates with larger MSPEs than those obtained from the filtering process, again indicating that the underlying statistical model is no longer valid. When we applied smoothing to the estimates obtained by applying the location error adjusted updating equations, there was again an improvement in the MSPEs. In summary, the two dimensional simulation results suggest that the adjusted updating equations performs better than the usual updating equations, when errors affect the measured location of system observations. This is in complete agreement with the one dimensional simulations.

## Discussion

Our motivation for this study was the challenge of utilizing the voluminous amounts of sensor data collected from marine animals for oceanographic data assimilation schemes. In this article we demonstrated how to make simple adjustments to a standard Kalman filtering DA scheme to account for the uncertainty in the spatial location of the observations. Our simulations demonstrated that that the adjustments were effective and gave better results than a standard Kalman filtering DA scheme, given error prone location estimates. While our example is simple, we note that our adjustment method can also be extended to non-linear filtering schemes like the extended Kalman filter (EKF; e.g., [20]) or the ensemble Kalman filter (EnKF; e.g., [3]).

The EnKF is a Monte-Carlo implementation of the Kalman filter algorithm and is often used to data assimilate oceanographic models (e.g., [3]). Location-error adjustments to the Kalman filter detailed here will also apply to the EnKF in a simple and straightforward way. In the EnKF algorithm, at each step, one needs to simulate observations $\mathbf{Y}_t$, by adding N(0,$\mathbf{v}_t$) noise, (recalling that $\mathbf{v}_t$ is the observation error) to build an ensemble of model predictions. As before, including location error inflates the variance of $\mathbf{Y}_t$. Hence, as demonstrated above, the correct variance when replicating $\mathbf{Y}_t$ is not N(0,$\mathbf{v}_t$) random errors but N(0, $V(\mathbf{Y}_t | \boldsymbol{\xi}_t, \mathbf{Y}_{\leq t-1})$) random errors, where $V(\mathbf{Y}_t | \boldsymbol{\xi}_t, \mathbf{Y}_{\leq t-1})$ is given by (8).

Throughout this article we assumed that both the system equation and the observation equation are linear in $\mathbf{Z}_t$'s. However, in many applications this is not the case, and the EKF was developed to tackle this particular problem. In the EKF, a linearized approximation, defined by the Jacobian, or the linear tangent operator, is used for the prediction of the error statistics. The algorithm is otherwise quite similar to the simple Kalman filter algorithm. Therefore, the adjustment for this case will be along the same lines as the corrections proposed above.

In this article, we have assumed that the process variable $\mathbf{Z}_t$ is independent and has a diagonal covariance matrix $\mathbf{W}_t$. This is a simplifying assumption that enables the Kalman filter to be fitted with relative ease. However, process variables are often spatially correlated at any given time step. This correlation could be incorporated into our altered filtering scheme by specifying the covariance matrix $\mathbf{W}_t$ to be an appropriate, non-diagonal, structure. The results obtained in this article will still hold with this alteration.

Another simplifying assumption that we made in this article is the independence of the location errors. However, in many operational situations, location errors might not be independent over time, instead they might be auto-correlated. The adjustments to the Kalman filter based DA scheme may need further refinement from those described here, depending on the precise factor giving rise to auto-correlated location error.

In this article we consider data assimilation schemes that use Kalman filter based methods. However, these methods are not the only ones used in practice. A notable alternative are the class of variational methods (e.g., [2]). In variational methods the location of the observations are still required. A description of how this extra uncertainty should be incorporated into a variational method remains a topic of future research.

The simulation study illustrated in this article showed that the performance of the adjusted Kalman filter updates is better than the unadjusted Kalman filter updates for data prone to location error. This study has charted a way forward to address this problem. However, the methods shown here require further development and expansion into a real ocean data assimilation

scheme, in order to assess the performance in the context of more complicated oceanographic models.

## Supporting Information

**Appendix S1  Glossary of notations.**
(PDF)

**Appendix S2  First-order corrections to the Kalman filter algorithm.**
(PDF)

**Appendix S3  Second-order correction to the posterior expectation.**
(PDF)

## Author Contributions

Conceived and designed the experiments: MB TP SF AS. Performed the experiments: AS SF. Contributed reagents/materials/analysis tools: AS TP SF. Wrote the paper: AS TP SF MB. Identified problem: TP MB. Generated details of the method and the study design: AS SF MB. Wrote the manuscript: TP MB AS SF.

## References

1. Oke PR, Schiller A (2007) Impact of Argo, SST, and alimeter data on an eddy-resolving ocean reanalysis. Geophysical Research Letters 34.
2. Ménard R, Daley R (1996) The application of Kalman smoother theory to the estimation of 4DVAR error statistics. Tellus A 48: 221–237.
3. Evensen G (2003) The ensemble Kalman filter: Theoretical formulation and practical implementation. Ocean Dynamics 53: 343–367.
4. Costa DP, Block BA, Bograd SJ, Fedak MA, Gunn JS (2010) TOPP as a marine life observatory: Using electronic tags to monitor the movements, behaviour and habitats of marine vertebrates. In: Hall J, Harrison DE, Stammer D, editors, In Proceedings of OceanObs'09: Sustained Ocean Observations and Information for Society (Vol. 1). Venice, Italy, pp. 21–25.
5. Niller PP, Paduan JD (1995) Wind-driven motions in the Northeast Pacific as measured by Lagrangian drifters. Journal of Physical Oceanography 25: 2819–2830.
6. Poulain PM, Warn-Varnas A, Niller PP (1996) Near-surface circulation of the Nordic seas as measured by Lagrangian drifters. Journal of Geophysical Research 101: 18237–18258.
7. Charrassin JB, Hindell M, Rintoul SR, Roquet F, Sokolov S, et al. (2008) Southern ocean frontal structure and sea-ice formation rates revealed by elephant seals. Proceedings of the National Academy of Sciences 105: 11634.
8. Boehme L, Meredith MP, Thorpe SE, Biuw M, Fedak M (2008) Antarctic Circumpolar Current frontal system in the South Atlantic: Monitoring using merged Argo and animal-borne sensor data. Journal of Geophysical Research 113: C09012.
9. Lydersen C, Nost OA, Lovell P, McConnell BJ, Gammelsrod T, et al. (2002) Salinity and temperature structure of a freezing Arctic fjord-monitored by white whales (Delphinapterus leucas). Geophysical Research Letters 29: 34–1.
10. Ekstrom PA (2004) An advance in geolocation by light. Mem Nat Instit Polar Res, Special Issue: 210–226.
11. Gunn J, Block B (2001) Advances in acoustic, archival, and satellite tagging of tunas. In: Tuna: Physiology, Ecology, and Evolution, Academic Press, volume 19 of Fish Physiology. pp. 167–224.
12. Metcalfe JD, Arnold GP (1997) Tracking fish with electronic tags. Nature 387: 665–666.
13. Teo SL, Boustany A, Blackwell S, Walli A, Weng KC, et al. (2004) Validation of geolocation estimates based on light level and sea surface temperature from electronic tags. Marine Ecology Progress Series 283: 81–98.
14. Pedersen MW (2007) Hidden Markov models for geolocation of fish. Master's thesis, Informatics and Mathematical Modelling, Technical University of Denmark, DTU, Richard Petersens Plads, Building 321, DK-2800 Kgs. Lyngby. URL http://www2.imm.dtu.dk/pubdb/p.php?5183. [Last accessed online: 10 July, 2012].
15. Vincent C, Mcconnell BJ, Ridoux V, Fedak MA (2002) Assessment of Argos location accuracy from satellite tags deployed on captive gray seals. Marine Mammal Science 18: 156–166.
16. Royer F, Lutcavage M (2009) Positioning pelagic fish from sunrise and sunset times: Complex observation errors call for constrained, robust modeling, Springer Netherlands, volume 9 of Reviews: Methods and Technologies in Fish Biology and Fisheries. pp. 323–341.
17. Fedak M, Lovell P, McConnell B, Hunter C (2002) Overcoming the constraints of long range radio telemetry from animals: Getting more useful data from smaller packages. Integrative and Comparative Biology 42: 3–10.
18. Sibert JR, Nielsen A, Musyl MK, Leroy B, Evans K (2009) Removing bias in latitude estimated from solar irradiance time series. Tagging and Tracking of Marine Animals with Electronic Devices : 311–322.
19. Meinhold J, Singpurwalla ND (1983) Understanding the Kalman filter. The American Statistician 37: 123–127.
20. Anderson BD, Moore JB (1979) Optimal Filtering. Englewood Cliffs, NJ: Prentice-Hall, Inc.
21. Musyl MK, Brill RW, Curran DS, Gunn JS, Hartog JR, et al. (2001) Ability of archival tags to provide estimates of geographical position based on light intensity. In: Electronic tagging and tracking in marine fisheries: proceedings of the Symposium on Tagging and Tracking Marine Fish with Electronic Devices, February 7–11, 2000, East-West Center, University of Hawaii. Kluwer Academic Pub, volume 1, pp. 343–367.
22. Harvey AC (1989) Forecasting, Structural Time Series Models and the Kalman Filter. Cambridge, UK: Cambridge University Press.
23. Ide K, Courtier P, Ghil M, Lorenc AC (1997) Unified notation for data assimilation: Operational, sequential and variational. Journal of the Meteorological Society of Japan 75: 181–189.
24. Press WH, Teukolsky SA, Vetterling WT, Flannery BP (1988) Numerical Recipies in C++: The Art of Scientific Computing. Cambridge, UK: Cambridge University Press.
25. Nielsen A, Sibert J (2007) State-space model for light-based tracking of marine animals. Canadian Journal of Fisheries and Aquatic Sciences 64: 1055–1068.