

Point process modelling of the Afghan War Diary

Andrew Zammit-Mangion^{a,b}, Michael Dewar^c, Visakan Kadiramanathan^d, and Guido Sanguinetti^{a,e,1}

^aSchool of Informatics, University of Edinburgh, Edinburgh EH8 9AB, United Kingdom; ^bUniversity/British Heart Foundation Centre for Cardiovascular Science, Queen's Medical Research Institute, Edinburgh EH16 4TJ, United Kingdom; ^cDepartment of Applied Physics and Applied Mathematics, Columbia University, New York, NY 10027; ^dDepartment of Automatic Control and Systems Engineering, University of Sheffield, Sheffield S1 3JD, United Kingdom; and ^eSynthSys—Systems and Synthetic Biology, University of Edinburgh, Edinburgh EH9 3JD, United Kingdom

Edited by Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, and approved June 8, 2012 (received for review February 25, 2012)

Modern conflicts are characterized by an ever increasing use of information and sensing technology, resulting in vast amounts of high resolution data. Modelling and prediction of conflict, however, remain challenging tasks due to the heterogeneous and dynamic nature of the data typically available. Here we propose the use of dynamic spatiotemporal modelling tools for the identification of complex underlying processes in conflict, such as diffusion, relocation, heterogeneous escalation, and volatility. Using ideas from statistics, signal processing, and ecology, we provide a predictive framework able to assimilate data and give confidence estimates on the predictions. We demonstrate our methods on the WikiLeaks Afghan War Diary. Our results show that the approach allows deeper insights into conflict dynamics and allows a strikingly statistically accurate forward prediction of armed opposition group activity in 2010, based solely on data from previous years.

conflict prediction | point processes | variational Bayes

The last decade has witnessed a tremendous increase in the availability of data relating to conflicts. For example, the collection of media reports in the 'Armed Conflict Location and Event Dataset' (1) provides a small scale but highly curated record of conflict events. More prominently, the release of confidential documents by the WikiLeaks whistleblower website in July 2010 has provided for the first time a large scale (but uncurated) description of the current Afghan conflict. However, most analyses of these and similar data sources do not go beyond visualization and descriptive statistical methods (2–5), for good reasons: first, conflict data is highly heterogeneous and often poorly annotated. For example, the WikiLeaks Afghan War Diary (AWD) data used in this study (Dataset S1) consists of event entries as diverse as elaborate preplanned military activity and spontaneous stop-and-search events. Any plausible attempt to model this data will need to be statistical in nature in order to handle the high levels of noise. Second, it is very difficult to define simple mechanisms that would allow the bottom-up construction of a plausible model.

Here, we develop statistical dynamical modelling methodologies to provide a predictive framework that may be used in policy making. We show that the temporal and spatial dependencies (6, 7) as well as diffusion and advection effects (8, 9) inherent in conflict data make it suitable for the use of a broad class of models, widely employed in ecology and epidemiology, in order to describe the dynamics of disaggregated data. We then develop tools based on ideas from point process statistics (10) to constrain the models. The approach enables us to leverage powerful techniques from point process filtering theory and spatiotemporal statistics (11–14) to carry out inference of the underlying system's dynamics and to predict the future behavior of the system.

We test the performance of our methods on the AWD, a WikiLeaks release which contains over 75,000 military logs by the USA military, describing events which occurred between the beginning of 2004 and the end of 2009 and providing a high temporal and spatial resolution description of the Afghan war in that period. We show that our approach allows deeper insights in the conflict dynamics than simple descriptive methods by providing a spatially resolved map of the growth and volatility of the conflict.

Most remarkably, we show that a model trained on the AWD can predict with surprising statistical accuracy the progression of the conflict in 2010; i.e., a year after the end of the AWD data. We conclude the paper by discussing the importance and potential of statistical modelling of conflict data, as well as offering some consideration as to its wide applicability.

Statistical Methods

Spatial Point Processes and the Stochastic Integro-Difference Equation (SIDE). Conflict data typically consists of a set of incidents labeled through spatiotemporal coordinates which, when visualized as event markers, are highly spatiotemporally correlated, generally clustered and representative of some underlying structure. In this regard, these data sets are very similar to others encountered in a variety of fields, such as epidemiology (15) and agricultural sciences (16). Poisson point processes provide a convenient and frequently used mathematical framework to model event-based data; in this framework, the probability of observing a certain number of events within a region of interest \mathcal{O} is given by a Poisson distribution whose mean is the integral over \mathcal{O} of an intensity function $\lambda(s)$, $s \in \mathcal{O}$. In order to accommodate phenomena such as event clustering, the intensity itself is often modeled as a random function, giving rise to *doubly stochastic* or *Cox* processes. A popular class of Cox processes, which will also be considered here, is the log-Gaussian Cox process (LGCP) where the logarithm of the event intensity is assumed to be a Gaussian process (GP). We recall that a GP is wholly defined by (i) a mean function $\mu(s)$ describing a global trend and (ii) a covariance function $k(s, r)$ indicating how the field at distinct points in space (s and r) covary (17).

Because conflict data is often logged in a discrete-time format (e.g., the day of an event as opposed to the precise time), we will consider a discrete-time series of continuous-space LGCPs. Formally, let $k \in \mathcal{K}$, $\mathcal{K} = \{1, 2, \dots, K\}$ denote a discrete-time index set and $\{z_k(s)\}$, $z_k(s) \sim \mathcal{GP}(\mu_k(s), \sigma_k^2 \Psi_k(s, r))$, a set of temporally correlated spatial GPs, each with mean $\mu_k(s)$ and covariance function $\sigma_k^2 \Psi_k(s, r)$. For each k , we then define the point process intensity function as $\lambda_k(s) = \exp(z_k(s))$. Frequently, the mean function of $z_k(s)$, $k \in \mathcal{K}$, can be related to explanatory variables, such as population density, which help to reduce prediction uncertainty. We hence let $\mathbf{d}(s)$ be a vector of spatially referenced covariates and \mathbf{b}^T the corresponding regression parameters; the LGCP at time k then has intensity $\lambda_k(s) = \exp(\mathbf{b}^T \mathbf{d}(s) + z_k(s))$.

Naturally, the key question is how to specify the temporal dynamics of the intensity functions through $z_k(s)$; we need a sufficiently flexible modelling approach to incorporate the complexity of conflict dynamics. One such representation is the stochastic

Author contributions: A.Z.M., V.K., and G.S. designed research; A.Z.M. and G.S. performed research; A.Z.M., M.D., and G.S. analyzed data; and A.Z.M. and G.S. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence should be addressed. E-mail: G.Sanguinetti@ed.ac.uk.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1203177109/-DCSupplemental.

integro-difference equation (SIDE), a model originally introduced in ecology (18) which has rapidly gained popularity in spatiotemporal statistics (19). The SIDE relates the spatiotemporal dependent variable $z_k(\mathbf{s})$ to $z_{k+1}(\mathbf{s})$ through the following integral equation

$$z_{k+1}(\mathbf{s}) = \int_{\mathcal{O}} k_I(\mathbf{s}, \mathbf{r}) f_1(z_k(\mathbf{r})) d\mathbf{r} + e_k(\mathbf{s}), \quad [1]$$

where $k_I(\mathbf{s}, \mathbf{r})$ is the mixing kernel in the integral and $e_k(\mathbf{s})$ is an added disturbance, modeled as a Gaussian field with mean $\mu_Q(\mathbf{s})$ and covariance function $k_Q(\mathbf{s}, \mathbf{r})$, $e_k(\mathbf{s}) \sim \mathcal{GP}(\mu_Q(\mathbf{s}), k_Q(\mathbf{s}, \mathbf{r}))$, and \mathcal{O} is the spatial domain under investigation. The nonlinear mapping $f_1(\cdot)$ distorts the field in the *sedentary stage*; in this work we will employ the identity map $f_1(z_k(\mathbf{r})) = z_k(\mathbf{r})$, an assumption usually adopted in the absence of a priori knowledge (20). The SIDE is, in its original form, a very flexible modelling tool, capable of representing a number of dynamic effects such as diffusion and dispersal (or both simultaneously) even under considerably restrictive conditions (19). Although the AWD will suggest the use of only a special case of SIDE, the two-pronged methodological approach we present here to estimate unknown components is in principle applicable to the more general case.

Nonparametric Analysis. We start by studying the correlation between the conflict events within the same and across subsequent time frames. We are interested in the probabilities of finding a conflict event at \mathbf{r} given that an event has occurred at \mathbf{s} within the same time frame k or at the previous time frame $k-1$. In point process statistics these are quantified through the pair auto-correlation function (PACF) $g_{k,k}(\mathbf{s}, \mathbf{r})$, and what we term the pair cross-correlation function (PCCF) $g_{k,k+1}(\mathbf{s}, \mathbf{r})$ defined as

$$g_{k,k}(\mathbf{s}, \mathbf{r}) = \frac{\lambda_{k,k}^{(2)}(\mathbf{s}, \mathbf{r})}{\lambda_k^{(1)}(\mathbf{s})\lambda_k^{(1)}(\mathbf{r})}, \quad [2]$$

$$g_{k,k+1}(\mathbf{s}, \mathbf{r}) = \frac{\lambda_{k,k+1}^{(2)}(\mathbf{s}, \mathbf{r})}{\lambda_k^{(1)}(\mathbf{s})\lambda_{k+1}^{(1)}(\mathbf{r})}, \quad [3]$$

where $\lambda_k^{(1)}(\mathbf{s}) = \mathbb{E}[\lambda_k(\mathbf{s})]$ and $\lambda_{k,k}^{(2)}(\mathbf{s}, \mathbf{r}) = \mathbb{E}[\lambda_k(\mathbf{s})\lambda_k(\mathbf{r})]$ are real and positive and $\mathbb{E}[\cdot]$ denotes the expectation operator.

The PACF may be used to determine qualitative characteristics of the conflict; for instance if $g_{k,k}(\mathbf{s}, \mathbf{r}) = 1$, then no spatial pattern can be extracted from the data; $g_{k,k}(\mathbf{s}, \mathbf{r}) > 1$ and $g_{k,k}(\mathbf{s}, \mathbf{r}) < 1$ can be used to indicate conflict aggregation and repulsion respectively. The PACF can also be used as a preprocessing tool for dimensionality reduction. Direct use of the PACF and PCCF for nonparametric field estimation is also possible (SI Text) but our preliminary investigation showed that this is only a reliable proposition for homogeneous datasets with a very large number of events (SI Text).

Dimensionality Reduction and Bayesian Inference. In order to develop an inferential approach for SIDE driven LGCPs, we adopt a basis function representation of the spatiotemporal field, which we will then truncate at a level which enables sufficient accuracy (21). This representation, frequently employed in spatiotemporal modelling [e.g., process convolution models (22, 23)], in turn facilitates the implementation of computationally efficient inference algorithms.

The choice of basis functions is a problem that deserves attention; as far as we are aware, there are no standard solutions for LGCPs. We propose here a general approach to selecting basis functions based on the nonparametric estimation of the PACF. Specifically, we capitalize on (i) a fundamental lemma of LGCPs

$$g_{k,k}(\mathbf{s}, \mathbf{r}) = \exp(\sigma_k^2 \Psi_k(\mathbf{s}, \mathbf{r})), \quad [4]$$

which states that the log PACF is proportional to the field auto-correlation function and (ii) the auto-correlation theorem (24) which states that the Fourier transform of the auto-correlation function is the spectrum of the signal. Hence, a relationship between the frequency content of the point process and the PACF is found, which in turn may be used to select a set of sufficiently representative basis functions, much on the lines of refs. 21 and 25. We then obtain a decomposition of the kernel, the mean disturbance and the field as

$$z_k(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \mathbf{x}_k, \quad [5]$$

$$\mu_Q(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\theta}, \quad [6]$$

$$k_I(\mathbf{s}, \mathbf{r}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\Sigma}_I \boldsymbol{\phi}(\mathbf{r}), \quad [7]$$

$$k_Q(\mathbf{s}, \mathbf{r}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\Sigma}_Q \boldsymbol{\phi}(\mathbf{r}), \quad [8]$$

where $\boldsymbol{\phi}(\mathbf{s}) \in \mathbb{R}^n$ is the vector of basis functions, $\mathbf{x}_k \in \mathbb{R}^n$ and $\boldsymbol{\theta} \in \mathbb{R}^n$ are weights which reconstruct the spatiotemporal field and the disturbance mean respectively and where $\boldsymbol{\Sigma}_I \in \mathbb{R}^{n \times n}$ and $\boldsymbol{\Sigma}_Q \in \mathbb{R}^{n \times n}$ reconstruct the kernel covariance function and the disturbance covariance function respectively.

It can be shown (SI Text) that under this decomposition, the SIDE of Eq. 1 can be represented in the compact form

$$\mathbf{x}_{k+1} = \mathbf{A}(\boldsymbol{\Sigma}_I) \mathbf{x}_k + \mathbf{w}_k(\boldsymbol{\theta}, \boldsymbol{\Sigma}_Q), \quad [9]$$

where $\mathbf{A}(\boldsymbol{\Sigma}_I) \in \mathbb{R}^{n \times n}$ and $\mathbf{w}_k \in \mathbb{R}^n$ is a Gaussian colored noise term with mean $\mathbb{E}[\mathbf{w}_k] = \boldsymbol{\theta}$ and covariance $\text{cov}[\mathbf{w}_k] = \boldsymbol{\Sigma}_Q$. Eq. 9 is a standard linear dynamical system where both the states $\mathcal{X}_K = \mathbf{x}_{0:K} = \{\mathbf{x}_k\}_{k=0}^K$ and the unknown parameters $\boldsymbol{\theta} = \{\boldsymbol{\theta}, \boldsymbol{\Sigma}_I, \boldsymbol{\Sigma}_Q^{-1}\}$ need to be estimated from the data $\mathcal{Y}_K = \{\mathbf{y}_k\}_{k=1}^K$ where we define each \mathbf{y}_k to be the set of coordinates of the logged events at the k th time point.

For inference, we make use of the likelihood function

$$p(\mathbf{y}_k | \lambda_k(\mathbf{s})) = \prod_{s_j \in \mathbf{y}_k} \lambda_k(s_j) \exp\left(-\int_{\mathcal{O}} \lambda_k(\mathbf{s}) d\mathbf{s}\right), \quad [10]$$

and approximate each $\lambda_k(\mathbf{s})$ using the same basis representation:

$$\lambda_k(\mathbf{s}) = \exp(\mathbf{b}^T \mathbf{d}(\mathbf{s}) + z_k(\mathbf{s})) \approx \exp(\mathbf{b}^T \mathbf{d}(\mathbf{s}) + \boldsymbol{\phi}(\mathbf{s})^T \mathbf{x}_k). \quad [11]$$

We proceed with a computationally efficient variational Bayes (VB) method by approximating the full posterior distribution

$$p(\mathcal{X}_K, \boldsymbol{\theta}, \mathbf{b} | \mathcal{Y}_K) = p(\mathcal{X}_K, \boldsymbol{\theta}, \boldsymbol{\Sigma}_I, \boldsymbol{\Sigma}_Q^{-1}, \mathbf{b} | \mathcal{Y}_K) \approx \tilde{p}(\mathcal{X}_K) \tilde{p}(\boldsymbol{\theta}) \tilde{p}(\boldsymbol{\Sigma}_I) \tilde{p}(\boldsymbol{\Sigma}_Q^{-1}) \tilde{p}(\mathbf{b}), \quad [12]$$

where $\tilde{p}(\cdot)$ are the variational marginals (26, 27).

The variational marginals are able to reveal important properties of the conflict progression; \mathcal{X}_K is used to reconstruct the spatiotemporal field at every time point, $\boldsymbol{\theta}$ reveals the spatially varying escalation in conflict, $\boldsymbol{\Sigma}_I$ the extent of any spatial dynamics, if any, and $\boldsymbol{\Sigma}_Q$ the *volatility* of the conflict which can either be localized or dependent on events happening at remote geographical locations. The number of unknown parameters in the reduced model scales as $\mathcal{O}(n^2)$, where n is the number of basis functions retained. However, as we will see later, nonparametric

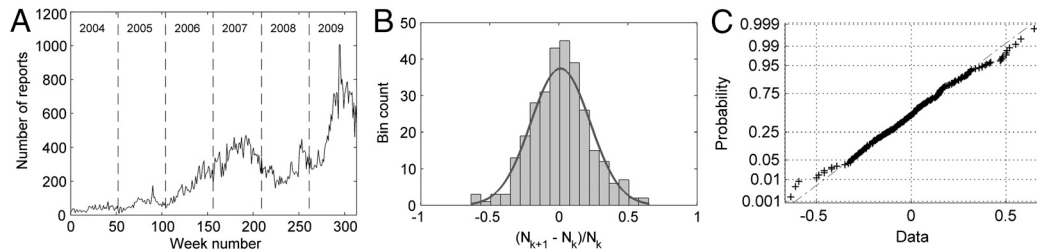


Fig. 1. Temporal analysis of the AWD. (A) Weekly number of activity reports in Afghanistan between January 2004 and December 2009 (bin size = 1 w). (B) Distribution of weekly fractional increments in report count in the AWD where N_k denotes the number of report counts at week k . (C) Corresponding normality probability plot. Fourteen points (4.5% of data) were marked as clear outliers as a result of low report count and not used in this analysis.

data analysis can suggest further simplifications which can considerably lower the complexity of the model.

The Afghan War Diary

On July 25, 2010, WikiLeaks publicly made available a compendium of US military war logs in Afghanistan dating between 2004 and 2009. The so-called Afghan War Diary contains a detailed insider's description of the military machinery of the world's largest power; it consists of roughly 77,000 logs and entries detail the time and position of an event, which could be anything from a stop-and-search episode to a gunfight. The dataset is considered a reliable description of the Afghan war and systematic verification efforts carried out by several organizations such as the New York Times* have found little reason to dispute its authenticity. *SI Text* reports some of our own tests which show significant correlations between the logged event rate in the AWD and that in other datasets. In what follows we adopt the spatiotemporal point process approach to infer a model from the data in the AWD and use it to analyze the heterogeneous growth (through ϑ) and volatility (through Σ_Q) of the conflict in Afghanistan and also to predict violence of armed opposition groups in 2010, a year after the end of the WikiLeaks dataset.

We start with a nonparametric analysis (*SI Text*) of the data by splitting the data into weekly intervals ($\Delta_t = 1$ week) and looking at the temporally averaged PACF and PCCF fitted to Gaussian radial basis functions. It is found that, on average, the log PACF is nearly identical to the log PCCF and that a nonparametric estimate of a homogeneous kernel $k_I(\|s - r\|)$, computed with the direct inverse filter, is very narrow in relation to the extent of the spatial correlations in the field (*SI Text*). This observation suggests that $k_I(\cdot)$ in the SIDE may be safely approximated to $\gamma(s)\delta(s - r)$, corresponding to negligible spatial interactions across adjacent time frames. Note that if $e_k(s)$ is restricted to be homogeneous and $\gamma(s) = \gamma$, the spatiotemporal covariance function is separable, a common assumption in several fields such as epidemiology (15). However, given the data characteristics, we chose to maintain the spatial heterogeneity in $e_k(s)$. We also set $\gamma(s) = 1$ as we found no evidence of mean reversion both at a national and a provincial level; additionally, we found that a spatially dependent $\gamma(s)$ did not contribute to increased prediction accuracy.

The resulting formulation is validated by studying the temporal dynamics of the AWD (Fig. 1A). A quantitative analysis reveals that the fractional increments of the event incidence nationwide are normally distributed (with a one-tailed Shapiro Wilk's test and a Levene's test with $\alpha = 0.1$, $n = 312$ w. See also Fig. 1B and C)[†]. This statistic characterizes systems following a geometric Brownian motion given by

$$d\lambda(s, t) = \tilde{R}(s)\lambda(s, t)dt + \lambda(s, t)dW(s, t), \quad [13]$$

where the increment $dW(s, t)$ is a Gaussian process with zero mean and covariance function $k_Q(s, r)dt$ and $\tilde{R}(s)$ is a spatially varying percentage drift. Applying Ito's Lemma (28) to $\ln \lambda(s, t)$ and noting that the continuous-time intensity $\ln \lambda(s, t) = b^T d(s) + z(s, t)$, we obtain the following form for $z(s, t)$:

$$dz(s, t) = R(s)dt + dW(s, t), \quad [14]$$

where $R(s) = \tilde{R}(s) - \frac{1}{2}\sigma(s)^2$ is a heterogeneous temporally independent spatial growth rate and $\sigma(s)^2$ is the variance field. Applying an explicit Euler discretization scheme to Eq. 14, one obtains the model $z_{k+1}(s) = z_k(s) + e_k(s)$ where $e_k(s)$ has mean $\mu_Q(s) = R(s)\Delta_t$ and covariance function $k_Q(s, r)\Delta_t$. This model is, as expected, the SIDE with the delta-Dirac kernel.

The field is next decomposed and Eqs. 6 and 8 are applied to finally obtain the random walk model occasionally employed in spatiotemporal studies (29)

$$x_{k+1} = x_k + w_k(\vartheta, \Sigma_Q). \quad [15]$$

For basis function selection we employed the aforementioned frequency-based approach (see *SI Text* for complete details). Finally, we chose population density and the distance to the nearest major city as covariates (see *SI Text* for details on how this choice was made). Inference was carried out using the VB algorithm described above. Full derivations, algorithmic details, and configuration parameters (priors and stopping conditions), as well as indicative run times, are given in *SI Text* respectively whilst a detailed simulation study showing the identifiability of the model under flat priors and a comparison with kernel-based estimators (30), is given in *SI Text*.

Results

Conflict Intensity and Regression Parameters. State inference leads to broad conclusions to where and how the conflict intensity has increased, decreased or shifted in time. We show the posterior mean intensity at regular intervals in *SI Text* and also in *Movie S1* together with the underlying AWD events at a weekly resolution. The progression of the intensity captures important geographical features of the war scenario. Regions of high intensity in 2009 include Sangin in northern Helmand (see *SI Text* for a provincial map), one of the most dangerous places in Afghanistan, notorious for thousands of improvised explosive devices and frequent suicide bombings (2). Other regions, such as Kabul, Nangarhar, and Paktya provinces, on the other hand have witnessed high activity throughout the six-year interval. Also very apparent is the emergence in later years of a high intensity ring starting from Kabul extending southwards towards Kandahar, up through Herat, through Balkh and back to Kabul. This roughly elliptical shape corresponds to the country's 'ring road', commonly targeted by insurgent activity and placement of improvised explosive devices (2). We note that a representative spatiotem-

*http://www.nytimes.com/2010/07/26/world/26editors-note.html?_r=1

[†]The Levene's test failed to reject the null hypothesis of constant variance for the years 2006 to 2009 but not when including 2004 and 2005. The reason for rejection when including the earlier two years can be safely attributed to relatively low report count, arising in noisy quantities when computing the fractional increments.

