Multi-task Gaussian process modelling for point process data

Virginia Aglietti

University of Warwick, The Alan Turing Institute

GP Summer School, September 2020

Outline

Motivation & Research Goal

- 2 Modelling point processes with GPs
- Crash course in multi-task GPs
 The MCPM model
- Crash course in variational inference for GPs
 Inference for the MCPM model
- 5 Experiments
 - Conclusions

Many social processes can be seen as point processes.

Many social processes can be seen as point processes.



- Cases of bovine tuberculosis (BTB) over the years 1989–2002 in Cornwall, UK.
- Four different strains of the disease.
- "Nonparametric estimation of spatial segregation in a multivariate point process: bovine tuberculosis in Cornwall, UK." by Diggle P. et al. (2005).

Many social processes can be seen as point processes.



• Counts of Covid19 events in UK



• Covid19 events in Tuscany

V.Aglietti

Many social processes can be seen as point processes.



- Cases of "assaults" over the year 2004-2013 in Chicago, USA.
- "Fast Kronecker Inference in Gaussian Processes with non-Gaussian Likelihoods" by Flaxman S. et al. (2015).

Many social processes can be seen as point processes.



Figure 1: Event counts for seven different crimes reported in 2016 in New York City, USA.

These social processes are characterized by:

- Spatial and temporal correlation;
- Correlation structure that changes in time and space (non-stationary) in a dependent manner (non-separable);
- Cross-correlation

Goal: Develop a scalable algorithm capable of jointly modelling social processes accounting for their complexities.

Cox processes

Estimating the intensity rate of events over a continuous space is a common problem for real-world applications.

Doubly stochastic Poisson process or **Cox process** [Cox 1955]:

- observed events are assumed to be generated from a Poisson process
- intensity function is modelled as another random process with a given prior probability measure

The ${\rm LGCP}\xspace$ model

The Log-Gaussian Cox Process (LGCP) [Møller et al., 1998] is an inhomogeneous Poisson process with a stochastic intensity function:

$$y_A|\lambda(\mathbf{x})\sim \mathsf{Poisson}\left(\int_{\mathbf{x}\in A}\lambda(\mathbf{x})d\mathbf{x}
ight)$$
 ,

where:

$$\lambda(\mathbf{x}) = \exp\{f(\mathbf{x})\}\$$

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}))$$

The Log-Gaussian Cox Process (LGCP) [Møller et al., 1998] is an inhomogeneous **Poisson process** with a **stochastic intensity** function:

$$y_A|\lambda(\mathbf{x})\sim \mathsf{Poisson}\left(\int_{\mathbf{x}\in\mathcal{A}}\lambda(\mathbf{x})d\mathbf{x}
ight)$$

The Log-Gaussian Cox Process (LGCP) [Møller et al., 1998] is an **inhomogeneous Poisson process** with a **stochastic intensity** function:

$$\gamma_A | \lambda(\mathbf{x}) \sim \mathsf{Poisson}\left(\int_{\mathbf{x} \in A} \lambda(\mathbf{x}) d\mathbf{x}
ight)$$

Computational grid for tractable inference:



The Log-Gaussian Cox Process (LGCP) [Møller et al., 1998] is an **inhomogeneous Poisson process** with a **stochastic intensity** function:

$$y_{\mathcal{A}}|\lambda(\mathbf{x})\sim \mathsf{Poisson}\left(\int_{\mathbf{x}\in\mathcal{A}}\lambda(\mathbf{x})d\mathbf{x}
ight)$$

Computational grid for tractable inference:

 $y_n | \lambda(\mathbf{x}) \sim \mathsf{Poisson}\left(\lambda(\mathbf{x}_n)\right)$

Limitations of the LGCP model

- Computational grid
- Single-task model
- I High computational cost due to the use of expensive MCMC schemes

Limitations of the LGCP model

- Omputational grid → "Structured variational inference in continuous cox process models." by Aglietti et al. (2019)
- Single-task model
- High computational cost when using expensive MCMC schemes

- Computational grid
- $\textbf{ igh computational cost} \rightarrow Variational Inference scheme$

Event counts for different types of crime are highly correlated.



Figure 2: Event counts for seven different crimes reported in 2016 in New York City, USA.

 \rightarrow Learn a model that capture the dependencies between these processes

Single-task GP regression

GP regression:

г



$$f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$$

$$y_i = f(\mathbf{x}_i) + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N \quad N = 3$$

$$\bigwedge \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ 0 & k_{13} \end{bmatrix} \right)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{bmatrix} + \sigma^2 I \right)$$

with $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j).$

Single-task GP regression



We can get the prediction for a new test point \mathbf{x}^* by computing $p(f(\mathbf{x}^*)|\mathbf{y})$ in closed form.

ł

Multi-task GP regression



$$f_{1}(\mathbf{x}) \sim \mathcal{GP}(0, k_{1}(\mathbf{x}, \mathbf{x}'))$$

$$f_{2}(\mathbf{x}) \sim \mathcal{GP}(0, k_{2}(\mathbf{x}, \mathbf{x}'))$$

$$\mathcal{D}_{1} = \{(\mathbf{x}_{i}^{1}, y_{i}^{1})\}_{i=1}^{N_{1}} \quad N_{1} = 3$$

$$\mathcal{D}_{2} = \{(\mathbf{x}_{i}^{2}, y_{i}^{2})\}_{i=1}^{N_{2}} \quad N_{2} = 2$$

$$y_{i}^{1} = f_{1}(\mathbf{x}_{i}^{1}) + \epsilon_{i} \quad \epsilon_{i} \sim \mathcal{N}(0, \sigma_{1}^{2})$$

$$y_{i}^{2} = f_{2}(\mathbf{x}_{i}^{2}) + \epsilon_{i} \quad \epsilon_{i} \sim \mathcal{N}(0, \sigma_{2}^{2})$$

$$\begin{bmatrix} \mathbf{y}^1 \\ \mathbf{y}^2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_2 \end{bmatrix} + \begin{bmatrix} \sigma_1^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \sigma_2^2 \mathbf{I} \end{bmatrix} \right)$$

Intrinsic coregionalization model (ICM)

We assume the two functions to be generated from latent process:

$$egin{aligned} u(\mathbf{x}) &\sim \mathcal{GP}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}')) \ f_1(\mathbf{x}) &= w_1 u(\mathbf{x}) \ f_2(\mathbf{x}) &= w_2 u(\mathbf{x}) \end{aligned}$$

$$\begin{bmatrix} \mathbf{y}^{1} \\ \mathbf{y}^{2} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \mathbf{K} + \begin{bmatrix} \sigma_{1}^{2} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \sigma_{2}^{2} \mathbf{I} \end{bmatrix} \right)$$
$$\mathbf{K} = \mathbf{B}k(\mathbf{x}, \mathbf{x}') = \mathbf{w}\mathbf{w}^{T}k(\mathbf{x}, \mathbf{x}') \quad \mathbf{w} = \begin{bmatrix} w_{1} & w_{2} \end{bmatrix}^{T}$$

Intrinsic coregionalization model (ICM)

In general we have a set of functions $\{f_{\rho}(\mathbf{x})\}_{\rho=1}^{P}$ and we assume

$$f_{p}(\mathbf{x}) = \sum_{i=1}^{R} w_{p}^{i} u^{i}(\mathbf{x})$$

with $\{u^i(\mathbf{x})\}_{i=1}^R$ are independent GPs with the same $k(\mathbf{x}, \mathbf{x}')$. Define $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), ..., f_P(\mathbf{x})]^T$ then we have:

$$\mathsf{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbf{B}k(\mathbf{x}, \mathbf{x}')$$

Other multi-task GP models

- Semiparametric Latent factor model (SLFM)
- Linear Coregionalization model (LCM)
- Process convolutions

• ...

A multi-task LGCP model

The **multivariate** LGCP (MLGCP) [Diggle et al., 2013] considers P types of points with an intensity given by:

$$\lambda_{\rho}(\mathbf{x}) = \exp(\beta + f_0(\mathbf{x}) + f_{\rho}(\mathbf{x}))$$

where $f_0(\mathbf{x})$ indicates a GP common to all tasks while $f_p(\mathbf{x})$ denotes a GP specific to the point process of type p. Inference proceeds via a MALA algorithm.

Coregionalisation models for point processes

The intensity function is a linear combination of Q independent Gaussian processes:

$$\lambda_p(\mathbf{x}) = \exp(\sum_{q=1}^{Q} \underbrace{w_{pq}}_{deterministic} u_q(\mathbf{x})).$$

The covariance $\text{Cov}[f_p(\mathbf{x}), f_{p'}(\mathbf{x}')]$ is given by $K(\mathbf{x}, \mathbf{x}') = \sum_{q=1}^{Q} B_q k_q(\mathbf{x}, \mathbf{x}')$ where B_q is known as a coregionalisation matrix.

A Bayesian treatment has been proposed by Schmidt & Gelfand (2003). The weight parameters have an inverse Wishart prior and inference proceeds with MCMC.

Modelling goals

- increase modelling flexibility with respect to MLGCP and ICM
- develop a full Bayesian model that propagates uncertainty in the mixing weights

Multi-task Cox process model (MCPM)

Latent functions:

• Q uncorrelated GP:

$$egin{aligned} p(\mathbf{F}|m{ heta}) &= \prod_{q=1}^Q p(\mathbf{F}_{ullet q}|m{ heta}_q) \ &= \prod_{q=1}^Q \mathcal{N}(\mathbf{F}_{ullet q}|m{0}, \mathbf{K}^q_{ ext{xx}}), \end{aligned}$$

where θ_q are the hyper-parameters for the *q*-th latent function.

Mixing weights:

• **Coupled weights**, independent across latent processes:

$$p(\mathbf{W}|\boldsymbol{ heta}_w) = \prod_{q=1}^{Q} \mathcal{N}(\mathbf{W}_{ullet q}; \mathbf{0}, \mathbf{K}_w^q),$$

where θ_w denotes the hyper-parameters.

Graphical model representation



Figure 3: MCPM- Plate diagram. $X'_{pd'}$ represents the inputs for the GP prior on **W**. When placing a Normal prior on each w_{pq} , we introduce the additional factorization across *P* (dashed plate).

Other models

- Sigmoidal transformation: Adams et al. (2009), Gunter et al. (2014), Lian et al. (2015), Donner and Opper (2018).
- Square transformation: Walder and Bishop (2017), Lloyd et al. (2015), Lian et al. (2015), Lloyd et al (2016), John and Hensman (2018).

Inference in GP modulated point process models

GP regression:

- Gaussian prior : $f(\mathbf{x}) \sim \mathcal{N}(0, k(\mathbf{x}, \mathbf{x}')) = p(f)$
- Gaussian likelihood: $y \sim \mathcal{N}(f, \sigma^2 I)$
- Gaussian posterior: $p(f|y) \propto \mathcal{N}(y|f, \sigma^2 I) \mathcal{N}(f|0, \sigma^2 I)$

PP model:

- Gaussian prior : $f(\mathbf{x}) \sim \mathcal{N}(0, k(\mathbf{x}, \mathbf{x}')) = p(f)$
- Non Gaussian likelihood:
 y ~ Poisson(exp(f))
- Non Gaussian posterior: $p(f|y) = \frac{p(f)(y|f)}{\int p(f)(y|f)df}$

Inference in GP modulated point process models

GP regression:

- Gaussian prior : $f(\mathbf{x}) \sim \mathcal{N}(0, k(\mathbf{x}, \mathbf{x}')) = p(f)$
- Gaussian likelihood: $y \sim \mathcal{N}(f, \sigma^2 I)$
- Gaussian posterior: $p(f|y) \propto \mathcal{N}(y|f, \sigma^2 I) \mathcal{N}(f|0, \sigma^2 I)$

PP model:

- Gaussian prior : $f(\mathbf{x}) \sim \mathcal{N}(0, k(\mathbf{x}, \mathbf{x}')) = p(f)$
- **Non** Gaussian likelihood: y ~ Poisson(exp(f))
- Non Gaussian posterior: $p(f|y) = \frac{p(f)(y|f)}{\int p(f)(y|f)df}$
- Sampling methods to obtain samples from the posterior
- Approximation of the posterior with something of known form \rightarrow Variational inference

Variational inference

Idea: Minimize a divergence between the true posterior distribution p(f|y) and an approximate posterior q(f) of known form e.g. $q(f) = \mathcal{N}(\mu, \mathbf{C})$

$$extsf{KL}(q(f)||p(f|y)) = \int q(f) extsf{log} rac{q(f)}{p(f|y)} df$$

Adjust the variational parameters of q(f) to minimize the KL divergence.

Variational inference

$$\mathcal{KL}(q(f)||p(f|y)) = -\underbrace{(\mathbb{E}_{q(f)}[\log(f, y)] - \mathbb{E}_{q(f)}[\log(q(f))])}_{\mathsf{ELBO}} + \mathsf{log}p(y)$$

Minimizing the KL divergence is equivalent to maximizing the ELBO which is a lower bound on the log-evidence.

Inference goals for the $\ensuremath{\operatorname{MCPM}}$ model

- Use a VI algorithm
- Reduce the computational cost, especially when P is large
- Avoid mixing issues of MCMC methods
- Get rid of further Monte Carlo approximations in the optimisation procedure

We introduce an **augmented prior over the latent functions** [Titsias, 2009 and Bonilla et al., 2016] defined by:

$$p(\mathbf{U}|\boldsymbol{\theta}) = \prod_{q=1}^{Q} \mathcal{N}(\mathbf{U}_{\bullet q}; \mathbf{0}, \mathbf{K}_{zz}^{q}) \text{ and } p(\mathbf{F}|\mathbf{U}, \boldsymbol{\theta}) = \prod_{q=1}^{Q} \mathcal{N}(\mathbf{F}_{\bullet q}; \tilde{\boldsymbol{\mu}}_{q}, \widetilde{\mathbf{K}}^{q}),$$

where $\tilde{\mu}_q = \mathbf{K}_{xz}^q (\mathbf{K}_{zz}^q)^{-1} \mathbf{U}_{\bullet q}$ and $\widetilde{\mathbf{K}}^q = \mathbf{K}_{xx}^q - \mathbf{K}_{xz}^q (\mathbf{K}_{zz}^q)^{-1} \mathbf{K}_{zx}^q$.

 $U_{\bullet q}$ denotes the inducing process for $F_{\bullet q}$ computed in the $M \times D$ matrix Z_q of *inducing inputs* ($M \ll N$).

Idea: Minimise the *distance* between the true posterior distribution and the approximate posterior distribution, the variational distribution.

$$\underbrace{p(\mathbf{F}, \mathbf{U}, \mathbf{W} | \mathcal{D})}_{\text{True posterior distribution}} \rightarrow \underbrace{q(\mathbf{F}, \mathbf{U}, \mathbf{W} | \boldsymbol{\nu})}_{\text{Variational distribution}} = p(\mathbf{F} | \mathbf{U}) \underbrace{q(\mathbf{U} | \boldsymbol{\nu}_u)}_{\text{(1)}} \underbrace{q(\mathbf{W} | \boldsymbol{\nu}_w)}_{\text{(2)}}$$

(1)
$$q(\mathbf{U}|\nu_u) = \prod_{q=1}^{Q} \mathcal{N}(\mathbf{U}_{\bullet q}; \mathbf{m}_q, \mathbf{S}_q)$$

(2) $q(\mathbf{W}|\nu_w) = \prod_{q=1}^{Q} \mathcal{N}(\mathbf{W}_{\bullet q}; \omega_q, \Omega_q)$
where $\nu_u = \{\mathbf{m}_q, \mathbf{S}_q\}$ are the variational parameters.

 $\label{eq:main_station} \begin{array}{l} \mbox{Minimization of the KL divergence} \rightarrow \mbox{Maximisation of the evidence lower bound} \\ \mbox{(ELBO):} \end{array}$

$$\mathcal{L}_{elbo}(\boldsymbol{\nu}) = \mathcal{L}_{kl}(\boldsymbol{\nu}) + \mathcal{L}_{ell}(\boldsymbol{\nu}),$$
 (1)

$$\mathcal{L}_{kl}(\boldsymbol{\nu}) = -\mathrm{KL}(q(\mathbf{F}, \mathbf{U}, \mathbf{W} | \boldsymbol{\nu}) \| p(\mathbf{F}, \mathbf{U}, \mathbf{W}))$$
(2)

$$= -\mathrm{KL}(q(\mathbf{U}|\boldsymbol{\nu}_u) \| p(\mathbf{U})) - \mathrm{KL}(q(\mathbf{W}|\boldsymbol{\nu}_w) \| p(\mathbf{W})), \tag{3}$$

$$\mathcal{L}_{ell}(\nu) = \underbrace{\mathbb{E}_{q(\mathbf{F}, \mathbf{U}, \mathbf{W} | \nu)}[\log p(\mathbf{Y} | \mathbf{F}, \mathbf{W})]}_{(4)}$$

Can we compute this expectation in closed form?

$$= -\sum_{n=1}^{N}\sum_{p=1}^{P} \exp(\phi_p) \mathbb{E}_{q(\mathbf{F}_{n\bullet})q(\mathbf{W}_{p\bullet})} \left(\exp\left(\mathbf{W}_{p\bullet}\mathbf{F}_{n\bullet}\right)\right)$$
(5)

$$+\sum_{n=1}^{N}\sum_{p=1}^{P}\sum_{q=1}^{Q}\left[y_{np}(\omega_{pq}\mu_{nq}+\phi_{p})-\log(y_{np}!)\right]$$
(6)

$$\mathbb{E}\left[\exp\left(t\mathbf{W}_{\rho\bullet}\mathbf{F}_{n\bullet}\right)\right] = \mathrm{MGF}_{\mathbf{W}_{\rho\bullet}\mathbf{F}_{n\bullet}}(t) \tag{7}$$

where $MGF_{\mathbf{W}_{p\bullet}\mathbf{F}_{n\bullet}}(t)$ denotes the moment generating function of $\mathbf{W}_{p\bullet}\mathbf{F}_{n\bullet}$ in t. The random variable $\mathbf{W}_{p\bullet}\mathbf{F}_{n\bullet}$ is the sum of products of independent Gaussians and its MGF is thus given by:

$$\mathsf{MGF}_{\mathbf{W}_{p\bullet}\mathbf{F}_{n\bullet}}(t) = \prod_{q=1}^{Q} \frac{\exp\left[\frac{t\gamma_{pq}\tilde{\mu}_{nq} + \frac{1}{2}(\tilde{\mu}_{nq}^{2}K_{w}^{qp} + \gamma_{pq}^{2}\widetilde{K}^{qn})t^{2}}{1 - t^{2}K_{w}^{qp}\widetilde{K}^{qn}}\right]}{\sqrt{1 - t^{2}K_{w}^{qp}\widetilde{K}^{qn}}},$$
(8)

where the expectation is computed with respect to the prior distribution of $\mathbf{W}_{p\bullet}$ and $\mathbf{F}_{n\bullet}$; γ_{pq} is the prior mean of w_{pq} ; \widetilde{K}^{qn} denotes the variance of f_{nq} ; and K_w^{qp} is the variance of w_{pq} .

Computational complexity

The time complexity of the algorithm is of order $\mathcal{O}(M^3)$.

- The KL-divergence includes distributions over M-dimensional variables and P-dimensional variables (M \gg P). Its computational complexity is thus independent of N.
- The ELL term decomposes as a sum of expectations over N. This enables the use of stochastic optimization techniques thus also making this term independent of N.

Experimental comparison

Baselines:

- Variational LGCP model
- Variational formulation of ICM with Poisson likelihood implemented in GPflow

Performance measures:

- Root mean square error (RMSE)
- Negative log predicted likelihood (NLPL)
- Empirical coverage of the posterior counts distribution (EC)

Setting:

 We show transfer by partitioning the spatial extent in subregions and create missing data folds for each task.

Synthetic data experiment



Figure 4: Four related tasks evaluated at 200 evenly spaced points in [-1.5, 1.5]. The red intervals denote the 50 contiguous observations removed from the training set of each task.

	RMSE				NLPL				CPU
	1	2	3	4	1	2	3	4	time
MCPM-N	38.61	7.86	5.71	4.68	20.99	3.75	3.31	3.02	0.18
MCPM-GP	38.58	7.69	5.70	4.71	20.95	3.70	3.31	3.03	0.25
LGCP	48.17	14.32	11.83	5.38	43.40	8.78	8.98	3.27	0.32
ICM	39.07	7.96	7.88	6.03	21.81	3.76	3.77	3.38	0.52
		Empirica				erage (EC)		
		1		2		3		4	
MCPM-N	0.	80/0.1	2 0	.99 /0.	58 (0.92/0	.57	0.94/	0.83
MCPM-G	р О.	95/0.1	9 0.	72/ 0 .	67 1	.00/0	.78	0.92/	0.75
ICM	0.	75/0.0	3 0	.66/0.	60 (0.62/0	.50	0.93/	0.42

Table 1: Above: Performance on the missing intervals. MCPM-N and MCPM-GP denote independent and correlated prior respectively. Time in seconds per epoch. Lower values of NLPL are better. Below: In-sample/Out-of-sample 90% CI coverage for the predicted event counts distributions. Higher values of EC are better. Experiments

Transfer experiment - CRIME dataset



Figure 5: Estimated intensity surface when introducing missing data regions.

	Standardized NLPL (per cell)								
	1	2	3	4	5	6	7	time	
MCPM-N	0.56	0.91	0.66	1.09	0.85	10.29	0.42	2.05	
	(0.10)	(0.27)	(0.30)	(0.27)	(0.52)	(2.51)	(0.05)	2.85	
MCPM-GP	0.72	0.75	0.94	1.53	0.57	18.76	0.58	2 11	
	(0.18)	(0.18)	(0.55)	(0.52)	(0.19)	(8.25)	(0.12)	5.11	
LGCP	9.90	9.32	19.34	5.30	18.18	36.73	9.68	2.07	
	(3.66)	(2.41)	(11.45)	(1.02)	(8.65)	(4.02)	(2.67)	2.07	
ICM	0.87	1.36	0.91	1.19	0.69	12.30	0.93	44.12	
	(0.27)	(0.35)	(0.45)	(0.40)	(0.11)	(3.02)	(0.17)	44.15	

	Empirical Coverage (EC)							
	1	2	3	4	5	6	7	
MCPM+N	0.99/0.80	1.00/0.73	0.97/0.71	1.00/0.73	0.98/0.61	1.00/1.00	0.99/0.87	
MCPM-GP	1.00/0.87	1.00/0.74	1.00/0.71	1.00/0.95	1.00/0.88	0.80/1.00	1.00/0.85	
LGCP	0.86/0.29	0.76/0.20	0.86/0.29	0.82/0.37	0.68/0.25	0.94/0.00	0.83/0.21	
ICM	0.68/0.73	0.75/0.50	0.64/0.52	0.79/0.65	0.59/0.78	0.93/0.86	0.841/0.64	

Table 2: Right: CRIME dataset. Performance on the missing regions. Time in seconds per epoch. Lower values of NLPL are better. Left: In-sample/Out-of-sample 90% CI coverage for the predicted event counts distributions. Higher values of EC are better.

Transfer experiment - Bovine tuberculosis (BTB) dataset



Figure 6: First row: Counts of the BTB incidents on a 64×64 grid. Shaded areas represent missing data regions. Estimated conditional probabilities by MCPM (second row) and by ICM (third row).

Table 3: Upper: RMSE and NLPL on BTB with missing data. Time in seconds per epoch. Lower values of NLPL are better. Lower: In-sample/Out-of-sample 90% CI coverage for the predicted event counts distributions. Higher values of EC are better.

	RMSE				NLPL (per cell)				CPU
	gt 9	gt 12	gt 15	gt 20	gt 9	gt 12	gt 15	gt 20	time
MCPM-N	0.83	0.24	0.28	0.29	1.23	0.20	0.33	0.35	7.73
	(0.15)	(0.07)	(0.07)	(0.10)	(0.40)	(0.07)	(0.11)	(0.16))
MCPM-GP	(0.14)	(0.08)	(0.07)	(0.09)	(0.42)	(0.09)	(0.14)	(0.24)	7.63
LGCP	1.37	0.61	0.63	1.24	1.70	0.48	0.72	0.86	8 76
	(0.33)	(0.13)	(0.12)	(0.56)	(0.39)	(0.11)	(0.17)	(0.36)	
ICM	(0.15)	(0.07)	(0.08)	(5.48)	(0.40)	(0.06)	(0.10)	(0.14)	67.06
	Empirical Coverage (EC)								
		gt 9		gt 12	2	gt 1	5	gt 2	20
MCPM-N	0.	87/ 0.9)2 0	.97/ 0.	99	0.93/0	.96	0.95/1	1.00
MCPM-G	Р 0.	93 /0.9	91 0	.98/0.	98 (0.97/0	.98	0.97/	0.99
LGCP	0.	91/0.7	9 0	.97/0.	98	0.97/0	.97	0.96/0	0.98
ICM	0.	90/0.8	34 0	.96/0.	98	0.95/0	.96	0.96/0	0.96

Contributions:

- Multi-task GP models capture the complexity of different social processes
- Different multi-task GP models can be used to model the intensity function of a point process
- MCPM increases modeling flexibility wrt MLGCP and ICM by introducing stochastic mixing weights
- MCPM offers a fully Bayesian model that propagates uncertainty
- Posterior inference requires approximation when the likelihood is Poisson
- Scalable inference framework $(\mathcal{O}(M^3))$ considering prior and posterior distributions for **F** and **W** separately.

Efficient Inference in Multi-task Cox Process Models

Efficient Inference in Multi-task Cox Process Models

Virginia Aglietti University of Warwick The Alan Turing Institute Theodoros Damoulas University of Warwick The Alan Turing Institute Edwin V. Bonilla CSIRO's Data61 UNSW

Abstract

We generalize the log Gaussian Cox process (LGCP) framework to model multiple correlated point data jointly. The observations are treated as realizations of multiple LGCPs, whose log intensities are given by linear combinations of latent functions drawn from Gaussian process priors. The combination coefintensity determines event occurrences. Among these modeling approaches, the log Gaussian Cox process (LGCP, Møller et al., 1998) is one of the most wellestablished frameworks, where the intensity is driven by a Gaussian process prior (GP, Williams and Rasmussen, 2006). The flexibility of LGCP comes at the cost of incredibly hard inference challenges due to its doubly-stochastic nature and the notorious scalability issues of GP models. The computational problems are

What about more difficult correlation structure?

How can we write a multi-task GP model that incorporates the correlation structure existing in a causal graph?

Multi-task Causal Learnin	g with Gaussian Processes
---------------------------	---------------------------

Virginia Aglietti University of Warwick The Alan Turing Institute V.Aglietti@warwick.ac.uk

Mauricio Álvarez University of Sheffield Mauricio.Alvarez@sheffield.ac.uk Theodoros Damoulas University of Warwick The Alan Turing Institute T.Damoulas@warwick.ac.uk

Javier González Microsoft Research Cambridge Gonzalez.Javier@microsoft.com

Abstract

This paper studies the problem of learning the correlation structure of a set of intervention functions defined on the directed acyclic graph (DAG) of a causal model. This is useful when we are interested in jointly learning the causal effects of interventions on different subsets of variables in a DAG, which is common in field

Thanks for your attention!