

# Log-Gaussian Cox processes and line transect sampling: optimal design criteria

Kenneth Flagg<sup>a,\*</sup>, John Borkowski<sup>a</sup>, Andrew Hoegh<sup>a</sup>

<sup>a</sup>*Department of Mathematical Sciences, Montana State University, Bozeman, MT 59717*

---

## Abstract

*Goal of this paper (placeholder abstract).* Evaluate line-transect designs in terms of many design-based and model-based criteria for spatial prediction using Bayesian LGCP models. Identify promising criteria and initial designs for later optimization and sequential design (not actually optimizing yet). Illuminate any relationships among criteria that will be helpful for constrained optimization. Other innovations are to compute design-based criteria with respect to the path (i.e. as path integrals) and to introduce a model-based APV criterion with a variable weight that more strongly penalizes errors near a decision threshold.

*Keywords:* log-Gaussian Cox process, optimal sampling, model-based design, spatial sampling design

---

## 1. Introduction

Spatial point process models have long been considered generally infeasible because of their computational demands, but recent advances in Bayesian computing have made the Log-Gaussian Cox process an attainable model in practice (Rue et al., 2009; Lindgren et al., 2011; Illian et al., 2012; Simpson et al., 2016). Variable sampling effort leads to a degraded point pattern Chakraborty et al. (2011) and it is relatively simple to accomodate variable sampling effort in these models using modern computing tools (Yuan et al., 2017). However,

---

\*Corresponding author

Email address: `kenneth.flagg@montana.edu` (Kenneth Flagg)

the literature on optimal sampling for spatial point process models is in its  
10 infancy (Liu and Vanhatalo, 2020).

Point pattern data are routinely collected in species distribution studies and  
ordnance response projects. These applications may use quadrat sampling or  
line-transect sampling, with transect sampling being more common. When the  
objective is mapping where events occur in space, various spatial mapping proce-  
15 dures have been used. Traditionally these have involved aggregating the data to  
grid cell counts or computing moving averages. Aggregation has the downside of  
introducing arbitrary structure into the data by the choice of gridding scheme  
or averaging window, and requires unnecessary computation effort (Simpson  
et al., 2016). Software is now available to fit spatial point process models to  
20 data acquired via distance sampling and simultaneously estimate the detection  
function (Johnson et al., 2014; R Core Team, 2019).

In ecological settings, sampling plans are often designed around the goal of  
estimating total abundance. Ordnance response surveys are typically designed  
with the objective of detecting (but not necessarily mapping) intensity hotspots.  
25 However, to our knowledge, there has been very little work done in deciding  
*where* to collect data when the goal is to map the intensity using a spatial point  
process model. In this paper, we nearest-neighbor criteria to the path design  
setting, introduce a model-based criterion (TPAPV) that emphasizes precise  
prediction, and introduce a sequential path construction scheme as a starting  
30 point for future optimization. We then compare a variety of path design schemes  
with respect to a suite of model-based and design-based criteria for simulated  
point pattern data.

### 1.1. Spatial design

*Design-based sampling.* Most work done for points. (Or quadrats approximated  
35 as points?) Space-filling criteria may be good starting points (e.g nearest-  
neighbor distance). *add references about Latin hypercube sampling*

*Space-filling curves.* Used in circuit design (Fan et al., 2014) (*find more cita-  
tions*) and high-dimensional data visualization in bioinformatics (Anders, 2009).

Peano curve is very flexible for filling irregular shapes (Fan et al., 2014). Hilbert  
40 curve is easy to construct.

Space-filling curves are one-dimensional paths constructed iteratively; as the number of iterations goes to infinity, the limiting path has nonzero area and actually fills the space (Sagan, 1994). For applications we stop after a finite number of iterations.

45 *Model-based spatial design.* Regularity is optimal for spatial prediction but randomness and a variety of interpoint distances are best for parameter estimation (Diggle and Lophaven, 2006). Inhibitory plus close pairs is a good compromise (Chipeta et al., 2017).

### 1.2. Paths as sampling designs

50 While some ideas about the characteristics of a good point design apply to paths, creating an optimal path design is not as simple as connecting the points of a point design with line segments. There are many ways to connect points into a path, so optimal design criteria must apply to the whole path and not only to the waypoints.

55 Pollard et al. (2002) adaptively zigzagged their line transects in a species abundance survey.

The Visual Sample Plan software includes features to create systematic transect plans and augment plans with additional transects in regions lacking spatial coverage (Matzke et al., 2014). It helps the user choose the transect spacing to  
60 maximize the probability of detecting an intensity hotspot. However, it does not employ criteria to optimize spatial prediction.

Liu and Vanhatalo (2020) used narrow quadrats (swaths along line-transects) as their sampling units. The transects were short relative to the size of the study region and not connected into a path.

### 65 1.3. Multi-objective optimization

summarize Lark (2016) and related

we use a large suite of criteria to explore the relationships among them

## 2. Materials and methods

Heuristics of a good path design:

- 70 • Should start with a sparse design with regular spacing, then refine with infill
  - Provides good spatial coverage even if aborted early
  - Imagine downloading a high-resolution intensity jpeg over 56k
- Path should avoid sharp turns but is allowed to cross itself
- 75 • One option is to generate two segments at a time, first a short-to-medium length segment to get to the start of the next transect, then a medium-to-long segment for the transect
- Could have new segment length be negatively correlated with the previous segment length

### 80 2.1. Design-based Criteria

*Path length.* The total distance that traveled is often a constraint. Minimize it.

*Nearest neighbor distance.* A common criterion for space-filling designs, we adapt it to be meaningfully calculated for any point on a path. Define the  $k$ th-order nearest neighbor distance as  $\text{nnd}_k(u) = \min |u - v|$  where  $v$  is any point  
85 in the set of path segments at least  $k$  steps away from the segment containing  $u$ . If  $k = 0$ , this includes  $v$  in the same segment as  $u$  so trivially  $\text{nnd}_0(u) = 0$  for all  $u$  in the path.  $\text{nnd}_1(u)$  includes all segments except the one containing  $u$ .  $\text{nnd}_2(u)$  excludes the segment containing  $u$  and segments with which it shares vertices. Segments not accessible by a connected path starting at  $u$  are always  
90 included. Maximize  $\min[\text{nnd}_2(u)]$ ,  $\text{avg}[\text{nnd}_2(u)]$ , and  $\text{avg}[\text{nnd}_1(u)]$ .

### 2.2. Model-based Criteria

*Posterior prediction variance.* Minimize maximum prediction variance and average prediction variance for GP.

*Posterior parameter variance.* Minimize posterior variance for each model pa-  
 95 rameter (intercept, variance, spatial scale).

*Decision-based criteria.* Maximize error rates and AUC of thresholding the in-  
 tensity at an action level. Minimize the threshold-penalized average predictive  
 variance (TPAPV),

$$\text{TPAPV} = \int_{\mathcal{R}} \text{Var}[\lambda(u)] p^{|\lambda(u)-A|} du,$$

where  $A$  is the action level/decision threshold and  $0 < p < 1$  penalizes uncer-  
 tainty about the boundary in used for thresholding.

### 2.3. Sampling Schemes

*Parallel transects.* Parallel transects running the length of the site in the vertical  
 100 axis. Three ways of choosing the horizontal coordinate: simple random sample  
 (SRS), systematic with a random starting point and even spacing, inhibitory  
 plus close pairs.

*Latin hypercube sampling.* Random Latin hypercube design connected by short-  
 est path. Waypoints generated by the `lhs` R package Carnell (2020). Connected  
 105 into a the shortest path by the `TSP` package (Hahsler and Hornik, 2020).

*Space-filling curves.* Hilbert curve generated by `HilbertVis` package (Anders,  
 2009). This is a deterministic design, so a random offset is added.

*Particle movement model.* Models the way data are actually collected. Way-  
 points generated sequentially by generating a jump distance and a direction.  
 110 The jump distance is generated from a scaled beta distribution, and negatively  
 correlated with previous jump distance. This behavior should approximately  
 alternate between a short “transition” and a long transect. The negative corre-  
 lation was achieved by applying a  $1 - x$  transformation to a beta autoregressive  
 process (McKenzie, 1985). The direction angle is drawn from a bimodal dis-  
 115 tribution that is symmetric around 0 (a normal distribution reflectd about 0).  
*explain the Strauss part*

*set up to think about adaptive sampling (adding a transect at a time or stopping early but don't actually do it here)*

#### 2.4. Model fitting

120 INLA (Rue et al., 2009), SPDE (Lindgren et al., 2011), off-grid (Simpson et al., 2016)

#### 2.5. Simulation procedure

describe the simulation and site

### 3. Results

125 look at examples of designs that minimize each criterion  
look at examples of designs along the Pareto front

### 4. Discussion

discuss starting points for optimization and sequential design

### 5. Conclusions

## 130 Appendix A. Notation and Terminology

- process defined on  $\mathcal{D} \subset \mathbb{R}^d$ , domain of the intensity function, in this manuscript  $d = 2$
- observation window  $\mathcal{S} \subset \mathcal{D}$
- define three regions:
  - 135 – the domain  $\mathcal{D}$  over which the process mathematically operates
  - the study region  $\mathcal{R}$  over which inferences are desired
  - the observed/sampled observation window  $\mathcal{S}$

- general relationship is  $\mathcal{S} \subset \mathcal{R} \subset \mathcal{D} \subset \mathbb{R}^d$  where all of the subset symbols taken to mean “subset or equal”
- 140 •  $\mathcal{D}$  can be bounded or unbounded (often equal to  $\mathbb{R}^d$ ),  $\mathcal{S}$  practically always bounded,  $\mathcal{R}$  bounded or unbounded depending on application and inferential goals
- the “fully surveyed” (censused) situation is  $\mathcal{S} = \mathcal{R}$
- survey path  $\mathcal{P}$  is a one-dimensional subset of  $\mathcal{R}$ 
  - 145 – set of one or more sequences of waypoints connected by line segments
  - $\mathcal{S}$  is the set of all points within a fixed (and assumed known) radius of  $\mathcal{P}$
- $\mathbf{X}$  point process on  $\mathcal{R}$ ,  $\mathbf{x} = \{x_1, \dots, x_n\}$  realized point pattern
  - $\mathbf{X}_{\mathcal{S}} = \mathbf{X} \cap \mathcal{S}$  the restriction of  $\mathbf{X}$  to  $\mathcal{S}$ ,  $\mathbf{x} = \mathbf{X} \cap \mathcal{S}$  the realized
  - 150 observeable point pattern
- point  $x \in \mathbf{x}$  called an event
- intensity function  $\lambda(u)$
- types of “points” in space:
  - $x$  event in the point pattern
  - 155 –  $s$  numerical integration node
  - $u$  arbitrary location in  $\mathcal{D}$  used to index intensity function and predictors
- $z(u)$  a column vector of covariates/predictors at  $u$  (not used in this manuscript)
- “point” refers to a  $u$  unless clearly stated otherwise
- 160 • bold for sets and spatial processes, normal italics for spatial vectors
- $y$  and variations will be used for objects derived from the point pattern, e.g. marks, pseudodata

- distance sampling fits into the framework with expansion of notation to include a (nontrivial) detection function and differentiate between the observed and observable point patterns

165

## References

- Anders, S., 2009. Visualisation of genomic data with the Hilbert curve. *Bioinformatics* doi:10.1093/bioinformatics/btp152.
- Carnell, R., 2020. lhs: Latin Hypercube Samples. URL: <https://CRAN.R-project.org/package=lhs>. R package version 1.0.2.
- 170 Chakraborty, A., Gelfand, A.E., Wilson, A.M., Latimer, A.M., Silander, J.A., 2011. Point pattern modelling for degraded presence-only data over large regions. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 60, 757–776.
- 175 Chipeta, M., Terlouw, D., Phiri, K., Diggle, P., 2017. Inhibitory geostatistical designs for spatial prediction taking account of uncertain covariance structure. *Environmetrics* 28.
- Diggle, P., Lophaven, S., 2006. Bayesian geostatistical design. *Scandinavian Journal of Statistics* 33, 53–64.
- 180 Fan, J.A., Yeo, W.H., Su, Y., Hattori, Y., Lee, W., Jung, S.Y., Zhang, Y., Liu, Z., Cheng, H., Falgout, L., Bajema, M., Coleman, T., Gregoire, D., Larsen, R.J., Huang, Y., Rogers, J.A., 2014. Fractal design concepts for stretchable electronics. *Nature communications* 5, 3266.
- Hahsler, M., Hornik, K., 2020. TSP: Traveling Salesperson Problem (TSP). URL: <https://CRAN.R-project.org/package=TSP>. R package version 1.1-10.
- 185 Illian, J.B., Sørbye, S.H., Rue, H., 2012. A toolbox for fitting complex spatial point process models using integrated nested laplace approximation (inla). *The Annals of Applied Statistics* , 1499–1530.



- 190 Johnson, D., Laake, J., VerHoef, J., 2014. DSpat: Spatial Modelling for Distance Sampling Data. URL: <https://CRAN.R-project.org/package=DSpat>. R package version 0.1.6.
- Lark, R., 2016. Multi-objective optimization of spatial sampling. *Spatial Statistics* 18, 412–430.
- 195 Lindgren, F., Rue, H., Lindström, J., 2011. An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73, 423–498.
- Liu, J., Vanhatalo, J., 2020. Bayesian model based spatiotemporal survey designs and partially observed log Gaussian Cox process. *Spatial statistics* 35, 100392.
- 200 Matzke, B., Wilson, J., Newburn, L., Dowson, S., Hathaway, J., Sego, L., Bramer, L., Pulsipher, B., 2014. Visual Sample Plan Version 7.0 User’s Guide. Pacific Northwest National Laboratory. Richland, Washington. URL: <http://vsp.pnnl.gov/docs/PNNL-23211.pdf>.
- 205 McKenzie, E., 1985. An autoregressive process for beta random variables. *Management Science* 31, 988–997.
- Pollard, J., Palka, D., Buckland, S., 2002. Adaptive line transect sampling. *Biometrics* 58, 862–870.
- 210 R Core Team, 2019. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. Vienna, Austria. URL: <https://www.R-project.org/>.
- Rue, H., Martino, S., Chopin, N., 2009. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the royal statistical society: Series b (statistical methodology)* 71, 319–392.
- 215

Sagan, H., 1994. Space-filling curves. Springer.

Simpson, D., Illian, J.B., Lindgren, F., Sørbye, S.H., Rue, H., 2016. Going off grid: Computationally efficient inference for log-Gaussian Cox processes. *Biometrika* 103, 49–70.

220

Yuan, Y., Bachl, F.E., Lindgren, F., Borchers, D.L., Illian, J.B., Buckland, S.T., Rue, H., Gerrodette, T., et al., 2017. Point process models for spatio-temporal distance sampling data from a large-scale survey of blue whales. *The Annals of Applied Statistics* 11, 2270–2297.