

Recent Advances in Bayesian Spatial Survival Modeling

Tim Hanson

Department of Statistics
University of South Carolina, U.S.A.

University of Alabama
Department of Mathematics

April 10, 2015

Outline

- 1 Fundamental concepts
- 2 Semiparametric models
- 3 Spatial copula models

Survival data

- Can be time to any event of interest, e.g. death, leukemia remission, bankruptcy, electrical component failure, etc.
- Data T_1, T_2, \dots, T_n live in \mathbb{R}^+ .
- Called: survival data, reliability data, time to event data.
- Interest often focuses on relating aspects of the distribution on T_i to covariates or risk factors \mathbf{x}_i , possibly time-dependent $\mathbf{x}_i(t)$. Can be external or internal.

Survival data: covariates and censoring

- Uncensored data: $(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_n, t_n)$. Observe $T_i = t_i$.
- Right censored data: $(\mathbf{x}_1, t_1, \delta_1), \dots, (\mathbf{x}_n, t_n, \delta_n)$. Observe

$$\left\{ \begin{array}{ll} T_i = t_i & \delta_i = 1 \\ T_i > t_i & \delta_i = 0 \end{array} \right\}.$$

- Interval censored data: $(\mathbf{x}_1, a_1, b_1), \dots, (\mathbf{x}_n, a_n, b_n)$. Observe $T_i \in [a_i, b_i]$.

Density and survival

- Continuous T has density $f(t)$.
- Survival function is

$$S(t) = 1 - F(t) = P(T > t) = \int_t^{\infty} f(s) ds.$$

- Regression model: proportional odds.
- Probability of making it past 40 years is $S(40)$.
- Odds of dying before 40 years are $\frac{1-S(40)}{S(40)}$.

Quantiles

- p^{th} quantile q_p for T solves $P(T \leq q_p) = p$.
- Continuous $T \Rightarrow q_p = F^{-1}(p)$.
- Median lifetime in the population is $q_{0.5} = F^{-1}(0.5)$.
- Regression model: accelerated failure time (proportional quantiles).
- Quantile regression active area of research from frequentist & Bayesian perspectives, e.g. Koenker's excellent `quantreg` package for R.

Residual life

- Mean residual life

$$m(t) = E\{T - t | T > t\} = \frac{\int_t^\infty S(s) ds}{S(t)}.$$

- Can expect to live $m(40)$ more years given made it to 40.
- Regression model: proportional mean residual life.
- Also median residual life.

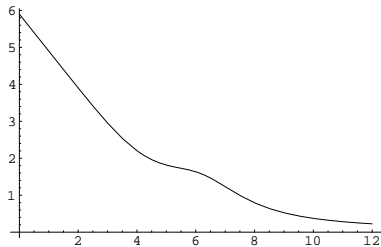
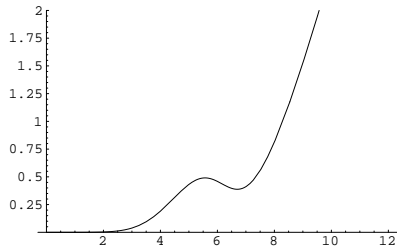
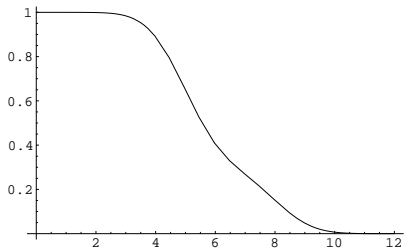
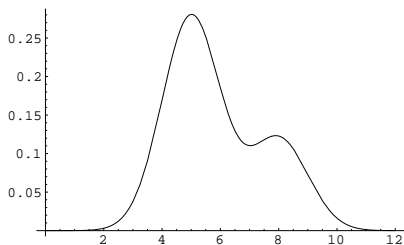
Hazard function

- Hazard at t :

$$h(t) = \lim_{dt \rightarrow 0^+} \frac{P(t \leq T < t + dt | T \geq t)}{dt} = \frac{f(t)}{S(t)}.$$

- Probability of dying tomorrow is $h(40) \left(\frac{1}{365}\right)$ given made it to 40 years.
- Regression models: proportional hazards (Cox), additive hazards (Aalen), accelerated hazards, & extended hazards.

Density, survival, hazard, and MRL



Nonparametric survival priors

- Infinite-dimensional process directly defined on one of $h(t)$, $H(t)$, $f(t)$, or $S(t)$.
- Note that prior on one function implies prior on other three.
- Priors on $h(t)$ include extended gamma, piecewise exponential, B-splines, etc.
- Priors on $H(t) = -\log S(t)$ include gamma, beta, etc.
- Priors on $S(t)$ include Dirichlet process (DP).
- Priors on $f(t)$ include DP mixtures, transformed Bernstein polynomials, Polya trees, B-splines, etc.
- Will consider B-spline and DP mixture.

Semiparametric models

Let's work covariates \mathbf{x}_i into the model for T_i . The most common way to do this is through a semiparametric model.

Why semiparametric?

- Splits inference into two pieces: β and $S_0(t)$ (or $h_0(t)$ or $m_0(t)$ or $H_0(t)$). Zero subscript stands for “baseline” where $\mathbf{x} = \mathbf{0}$.
- $\beta = (\beta_1, \dots, \beta_p)'$ succinctly summarizes effects of risk factors \mathbf{x} on aspects of survival.
- Make $S_0(t)$ as flexible as possible.
- Can make easily digestible statements concerning the population, e.g. “Median life on those receiving treatment A is 1.7 times those receiving B, adjusting for other factors.”

Some semiparametric models

- PH: $h_{\mathbf{x}}(t) = \exp(\mathbf{x}'\beta)h_0(t)$.
- AddH: $h_{\mathbf{x}}(t) = h_0(t) + \beta'\mathbf{x}$.
- AFT: $S_{\mathbf{x}}(t) = S_0\{e^{\beta'\mathbf{x}}t\}$.
- PO: $F_{\mathbf{x}}(t)/S_{\mathbf{x}}(t) = e^{\beta'\mathbf{x}}F_0(t)/S_0(t)$.
- PMRL: $m_{\mathbf{x}}(t) = e^{\beta'\mathbf{x}}m_0(t)$.
- AcchH: $h_{\mathbf{x}}(t) = h_0\{e^{\beta'\mathbf{x}}t\}$.
- ExtH: $h_{\mathbf{x}}(t) = h_0\{e^{\beta'\mathbf{x}}t\}e^{\gamma'\mathbf{x}}$.
- Others, but this covers 99%.

Proportional hazards (PH)

- Model is:

$$h_{\mathbf{x}}(t) = \exp(\mathbf{x}'\beta)h_0(t) \text{ or } S_{\mathbf{x}}(t) = S_0(t)^{\exp(\mathbf{x}'\beta)}.$$

- Stochastically orders $S_{\mathbf{x}_1}$ and $S_{\mathbf{x}_2}$.
- e^{β_j} is how risk changes when x_j is increased by unity.
- `BayesX` assigns penalized B-spline prior on $\log h_0(t)$ and allows for additive predictors, structured frailties, time-varying coefficients, etc. Free: <http://www.statistik.lmu.de/~bayesx/bayesx.html>. Also R package to call `BayesX`.
- Add `BAYES` command in `SAS PROC PHREG` gives p.w. exponential.
- Haiming Zhou's `spBayesSurv` has S_0 modeled as MPT.

Accelerated failure time (AFT)

- Model is

$$S_{\mathbf{x}}(t) = S_0 \left(e^{-\mathbf{x}'\beta t} \right), \quad \text{or} \quad \log T_{\mathbf{x}} = \mathbf{x}'\beta + e_0.$$

- Implies $q_p(\mathbf{x}) = e^{\mathbf{x}'\beta} q_p(0)$.
- Stochastically orders $S_{\mathbf{x}_1}$ and $S_{\mathbf{x}_2}$.
- e^{β_j} how any quantile – or mean – changes when increasing x_j by unity.
- Komarek's `bayesSurv` for AFT models; spline and discrete normal mixture on error.
- `bj()` in Harrell's `Design` library fits Buckley-James version.
- Haiming Zhou's `spBayesSurv` has S_0 modeled as MPT.

Proportional odds (PO)

- Model is

$$\frac{1 - S_{\mathbf{x}}(t)}{S_{\mathbf{x}}(t)} = \exp(\mathbf{x}'\beta) \frac{1 - S_0(t)}{S_0(t)}.$$

- e^{β_j} how odds of event occurring before t changes when x_j increased by unity (for any t).
- Attenuation of risk:

$$\lim_{t \rightarrow \infty} \frac{h_{\mathbf{x}_1}(t)}{h_{\mathbf{x}_2}(t)} = 1.$$

Plausible in many situations.

- Haiming Zhou's `spBayesSurv` has S_0 modeled as MPT.
`timereg` has frequentist version.

Spatial frailty survival models

- Survival data often collected over region.
- Georeferenced includes $\mathbf{s}_i = (x_i, y_i)$, e.g. latitude & longitude.
- Areal includes $c_i \in \{1, \dots, C\}$, e.g. the county of residence (there are C counties).
- Traditionally, spatial dependence induced by adding frailty (random effect) to linear predictor in semiparametric model.

Spatial frailty survival models

Georeferenced:

- Replace $\mathbf{x}'_i\beta$ by $\mathbf{x}'_i\beta + g_i$.
- Take $g_i = g(x_i, y_i)$ where $\{g(\mathbf{s}) : \mathbf{s} \in \mathcal{S}\}$ is mean-zero stationary Gaussian process.
- Yields $\mathbf{g} = (g_1, \dots, g_n) \sim N_n(\mathbf{0}, \mathbf{C}_\theta)$; \mathbf{C}_θ e.g. Matérn.

Areal:

- Replace $\mathbf{x}'_i\beta$ by $\mathbf{x}'_i\beta + g_{c_i}$.
- Define \mathbf{W} to be adjacency matrix: $w_{ij} = 1$ if counties i and j share a border, otherwise $w_{ij} = 0$ (assume $w_{ii} = 0$).
- CAR model assumes $g_j | \mathbf{g}_{-j} \sim N(\rho \tilde{g}_j, \frac{\lambda}{w_{j+}})$ where $\rho \in (0, 1)$ and $\tilde{g}_j = \frac{1}{w_{j+}} \sum_{i=1}^C w_{ij} g_i$.
- Limiting case $\rho \rightarrow 1$ called ICAR, requires $\sum_{j=1}^C g_j = 0$.

Both approaches provide mean-zero, smoothed spatial surface $g(\mathbf{s})$ or g_j over \mathcal{S} .

Spatial copula in a nutshell

- Let $T_i \sim F_{\mathbf{x}_i}(\cdot)$ where $F_{\mathbf{x}}$ c.d.f. from any survival model: parametric, semiparametric, nonparametric.
- $U_i = F_{\mathbf{x}_i}(T_i) \sim U(0, 1)$ and $Y_i = \Phi^{-1}(U_i) \sim N(0, 1)$. Let $\mathbf{Y} = (Y_1, \dots, Y_n)'$.
- No spatial correlation $\Rightarrow \mathbf{Y} \sim N_n(\mathbf{0}, \mathbf{I}_n)$.
- Spatial correlation $\Rightarrow \mathbf{Y} \sim N_n(\mathbf{0}, \mathbf{\Gamma})$. Here $\mathbf{\Gamma}_{n \times n} = [\gamma_{ij}]$ with pairwise correlations γ_{ij} .
- Li and Lin (2006) use this in PH model, term it “normal transformation model.”
- Gives marginal (population-averaged) model.
- Unlike frailties, can be used in models *without* a linear predictor.

SCCCR data set on prostate cancer survival

- Large dataset on prostate cancer survival that does not follow proportional hazards.
- $n = 20599$ patients from South Carolina Central Cancer Registry (SCCCR) for the period 1996–2004; each recorded with county, race, marital status, grade of tumor, and SEER summary stage; 72.3% are censored.
- Need to allow for non-proportional hazards and accommodate correlation of survival times within county.

Extended hazards model

- Etezadi-Amoli and Ciampi (1987) propose ExtH model

$$h_{\mathbf{x}}(t) = h_0(te^{\mathbf{x}'\beta})e^{\mathbf{x}'\gamma}.$$

- Say $\mathbf{x} = (x_1, x_2)$, then ExtH is

$$h_{\mathbf{x}}(t) = h_0(te^{\beta_1 x_1 + \beta_2 x_2})e^{\gamma_1 x_1 + \gamma_2 x_2}.$$

- $\gamma_1 = \beta_1 \Rightarrow x_1$ has AFT interpretation; $\beta_1 = 0 \Rightarrow x_1$ has PH interpretation; $\gamma_1 = 0 \Rightarrow x_1$ has Acch interpretation.

Baseline hazard $h_0(t)$

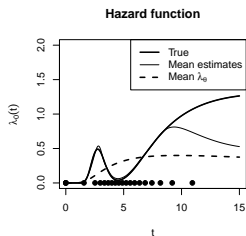
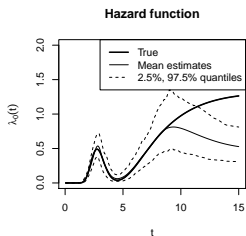
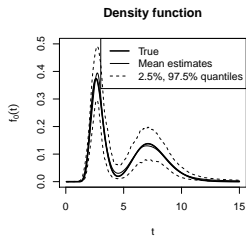
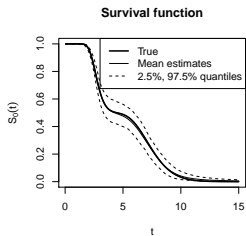
- Want to shrink $h_0(t)$ toward parametric target $h_\theta(t)$:

$$h_0(t) = \sum_{j=1}^J b_j B_{kj}(t)$$

where $B_{k1}(\cdot), \dots, B_{k,J}(\cdot)$ are k th order B-spline basis function over knots (s_1, \dots, s_{J+k}) .

- Let $\tilde{s}_j = \sum_{k=j}^{j+k} s_k / (k - 1)$ and $b_j = h_\theta(\tilde{s}_j)$.
- Schoenberg's approximation theorem (Marsden 1972) says $\max_{0 \leq t \leq s_{J+1}} \|h_0(t) - h_\theta(t)\| \leq \epsilon(h_\theta, k, J)$.
- Posterior updating: efficient MCMC with clever data augmentation.

Works great on simulated data



Spatial dependence via frailties impractical

- PH with frailties:

$$h(t_i|\mathbf{x}) = h_0(t_i)e^{\gamma'\mathbf{x}_i+g_{c_i}},$$

where g_{c_i} are county-level frailties, c_i is county subject i in.

- EH with frailties:

$$h(t_i|\mathbf{x}) = h_0\{t_i e^{\beta'\mathbf{x}_i+b_{c_i}}\}e^{\gamma'\mathbf{x}_i+g_{c_i}},$$

where, for our data, b_1, \dots, b_{46} and g_1, \dots, g_{46} are county-level frailties.

- Possible but impractical, and hard to interpret.

Spatial dependence via copula works great

- Define $Y_i = \Phi^{-1} \{F_{\mathbf{x}_i}(T_i)\}$.
- Under Li and Lin (2006) $\mathbf{Y} \sim N(\mathbf{0}, \Gamma)$.
- Likelihood from data $\{(t_i, \mathbf{x}_i, \delta_i)\}_{i=1}^n$ is

$$L(\beta, \gamma, \mathbf{b}, \theta, \Gamma) = \left[\prod_{i \in S} \frac{f_i(t_i)}{\phi(y_i)} \right] \int \left[\prod_{i \in S^c} \frac{f_i(z_i)}{\phi(y_i)} I(z_i > t_i) \right] \phi(\mathbf{y}; \mathbf{0}, \Gamma) \prod_{i \in S^c} dz_i$$

- How to define Γ ?
- We consider county-level lattice data; popular correlation model is intrinsic conditional autoregressive (ICAR) prior.

ICAR definition

- $\alpha = (\alpha_1, \dots, \alpha_C)$ is vector of correlated effects on over counties in S .
- ICAR prior on α is $p(\alpha) \propto \exp\{-\varphi \alpha'(\mathbf{D} - \mathbf{W})\alpha/2\}$
- Recall ICAR prior:

$$\alpha_j | \alpha_{-j}, \varphi \sim N\left(\sum_{j=1}^C w_{ij} \alpha_j / w_{j+}, 1/(\varphi w_{j+})\right).$$
- Random effects approach
 - $\tilde{\mathbf{Y}} = (\tilde{\mathbf{Y}}_1, \dots, \tilde{\mathbf{Y}}_C) = (\tilde{Y}_{11}, \dots, \tilde{Y}_{1n_1}, \dots, \tilde{Y}_{C1}, \dots, \tilde{Y}_{Cn_C})$.
 - $\tilde{Y}_{ij} = \alpha_i + \epsilon_{ij}$, $\alpha \sim N_C(\mathbf{0}, \Omega)$, $\epsilon \sim N_n(\mathbf{0}, \mathbf{I}\sigma^2)$.
 - Resulting correlation matrix $\Gamma = \text{corr}(\tilde{\mathbf{Y}})$ involves one unknown parameter φ^* .

Efficient evaluation of $\mathbf{y}'\Gamma^{-1}\mathbf{y}$

- Γ is a $n \times n$ matrix; needs to be inverted during MCMC.
- Elements of Γ^{-1} can be easily computed using SVD,

$$\Gamma^{-1} = \mathbf{A}^{-1}\mathbf{U}_1'((\mathbf{K}^* + \sigma^2\mathbf{I}_C)^{-1} - \sigma^{-2}\mathbf{I}_C)\mathbf{U}_1\mathbf{A}^{-1} + \sigma^{-2}\mathbf{A}^{-2}$$
 where \mathbf{A} is a diagonal matrix, $\mathbf{U}_1 = (\mathbf{u}_1, \dots, \mathbf{u}_C)$, \mathbf{u}_i is a vector of length n with ones corresponding to county i and zero elsewhere.
- $\mathbf{y}'\Gamma^{-1}\mathbf{y} = \mathbf{z}'((\mathbf{K}^* + \sigma^2\mathbf{I}_C)^{-1} - \sigma^{-2}\mathbf{I}_C)\mathbf{z} + \sigma^{-2}\mathbf{y}'\mathbf{A}^{-2}\mathbf{y}$ where $\mathbf{z} = \mathbf{U}_1'\mathbf{A}^{-1}\mathbf{y}$.

Savage-Dickey ratio for global and per-variable tests

- Example of global test of PH vs. EH

$$BF_{12} = \frac{\pi(\beta = \mathbf{0} | \mathcal{D}, EH)}{\pi(\beta = \mathbf{0} | EH)}.$$

- Example of per-variable of PH for x_j vs. EH

$$BF_{12} = \frac{\pi(\beta_j = 0 | \mathcal{D}, EH)}{\pi(\beta_j = 0 | EH)}.$$

SCCCR data

- SCCCR prostate cancer data for the period 1996–2004.
- Baseline covariates are county of residence, age, race, marital status, grade of tumor differentiation, and SEER summary stage.
- $n = 20599$ patients in the dataset after excluding subjects with missing information.
- 72.3% of the survival times are right-censored.

Goal: assess racial disparity in prostate cancer survival, adjusting for the remaining risk factors and accounting for the county the subject lives in.

SCCCR data

Table: Summary characteristics of prostate cancer patients in SC from 1996-2004.

| Covariate | | <i>n</i> | Sample percentage |
|--------------------|---|----------|-------------------|
| Race | Black | 6483 | 0.32 |
| | White | 14116 | 0.68 |
| Marital status | Non-married | 4525 | 0.22 |
| | Married | 16074 | 0.78 |
| Grade | well or moderately differentiated | 15309 | 0.74 |
| | poorly differentiated or undifferentiated | 5290 | 0.26 |
| SEER summary stage | Localized or regional | 19792 | 0.96 |
| | Distant | 807 | 0.04 |

Non-spatial EH and reduced models

Table: Summary of fitting the extended hazard model EH, the reduced model, AFT, and PH; * indicates $LPML - 21000$ and $DIC - 42000$.

| Covar | | EH | Reduced | AFT | PH | PH+additive age |
|----------------|------------|--------------------|----------------------|------------------|-----------------|-----------------|
| | | | | $\beta = \gamma$ | $\beta = 0$ | $\beta = 0$ |
| Age | β_1 | 0.50(0.48,0.52) | 0.48(0.46,0.50) | 0.48(0.45,0.51) | - | - |
| | γ_1 | 0.45(0.42,0.49) | $\gamma_1 = \beta_1$ | - | 0.65(0.62,0.68) | - |
| Race | β_2 | 0.18(0.15,0.21) | 0.20(0.16,0.21) | 0.18(0.15,0.22) | - | - |
| | γ_2 | 0.18(0.12,0.24) | $\gamma_2 = \beta_2$ | - | 0.26(0.21,0.32) | 0.26(0.20,0.31) |
| Marital status | β_3 | -0.06(-0.11,-0.02) | -0.05(-0.09,-0.00) | 0.26(0.21,0.30) | - | - |
| | γ_3 | 0.35(0.29,0.40) | 0.33(0.28,0.40) | - | 0.33(0.27,0.39) | 0.31(0.26,0.37) |
| Grade | β_4 | 0.03(-0.02,0.08) | $\beta_4 = 0$ | 0.27(0.22,0.32) | - | - |
| | γ_4 | 0.36(0.29,0.41) | 0.37(0.31,0.43) | - | 0.38(0.32,0.44) | 0.37(0.33,0.43) |
| SEER stage | β_5 | 3.19(2.80,3.53) | 3.27(2.79,3.57) | 1.50(1.41,1.59) | - | - |
| | γ_5 | 1.02(0.83,1.20) | 1.00(0.82,1.19) | - | 1.56(1.47,1.64) | 1.57(1.19,1.65) |
| $LPML^*$ | | -161.0 | -162.0 | -206.5 | -242.5 | -231.9 |
| DIC^* | | 267.7 | 270.7 | 366.0 | 443.0 | 412.8 |

Non-spatial EH and reduced models

Table: Bayes factors for comparing EH to PH, AFT, and AH with and without spatial correlation.

| Covariate | EH | | | Spatial+EH | | |
|----------------|--------|--------|--------|------------|--------|--------|
| | PH | AFT | AH | PH | AFT | AH |
| Age | > 1000 | 0.08 | > 1000 | > 1000 | 0.01 | > 1000 |
| Race | > 1000 | 0.01 | > 1000 | > 1000 | < 0.01 | > 1000 |
| Marital status | 1.79 | > 1000 | > 1000 | 1.18 | > 1000 | > 1000 |
| Grade | 0.14 | > 1000 | > 1000 | 0.08 | > 1000 | > 1000 |
| SEER stage | > 1000 | > 1000 | > 1000 | > 1000 | > 1000 | > 1000 |

Spatial EH and reduced models

Table: Summary of spatial models; * indicates *LPML* – 21000 and *DIC* – 42000.

| Covariates | | Marginal EH | Marginal reduced | PH+ICAR+additive age $\beta = 0$ |
|----------------|------------|--------------------|----------------------|-------------------------------------|
| Age | β_1 | 0.50(0.47,0.52) | 0.47(0.46,0.49) | – |
| | γ_1 | 0.46(0.43,0.49) | $\gamma_1 = \beta_1$ | – |
| Race | β_2 | 0.18(0.15,0.21) | 0.20(0.17,0.22) | – |
| | γ_2 | 0.17(0.11,0.23) | $\gamma_2 = \beta_2$ | 0.24(0.18,0.30) |
| Marital status | β_3 | -0.06(-0.10,-0.02) | -0.02(-0.05,-0.00) | – |
| | γ_3 | 0.34(0.28,0.41) | 0.33(0.27,0.39) | 0.32(0.25,0.38) |
| Grade | β_4 | 0.03(-0.01,0.07) | $\beta_4 = 0$ | – |
| | γ_4 | 0.36(0.30,0.42) | 0.38(0.32,0.43) | 0.37(0.32,0.44) |
| SEER stage | β_5 | 3.16(2.86,3.34) | 2.77(2.72,2.82) | – |
| | γ_5 | 1.10(0.94,1.26) | 1.21(1.01,1.33) | 1.55(1.46,1.64) |
| φ^* | | 50.1(19.9,113.7) | 54.6(22.7,120.8) | 33.08(9.2,100.1) |
| <i>LPML</i> * | | -142.7 | -143.2 | -215.7 |
| <i>DIC</i> * | | 192.4 | 164.0 | 332.5 |

Spatial EH and reduced models

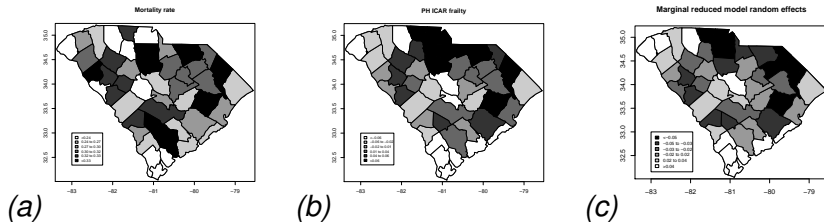


Figure: Map of (a) Mortality rate, (b) ICAR frailties in the PH model and (c) random effects in the marginal reduced model for SC counties.

Spatial EH and reduced models

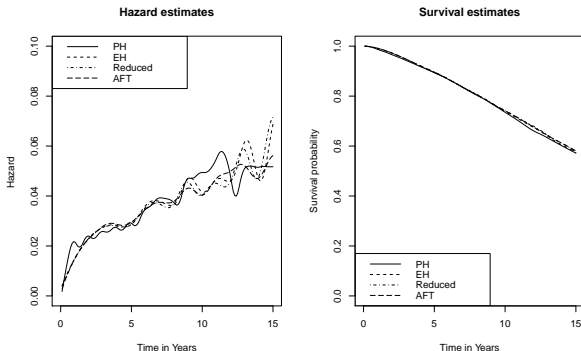


Figure: Baseline hazard (left) and survival probabilities (right) estimates.

Spatial EH and reduced models

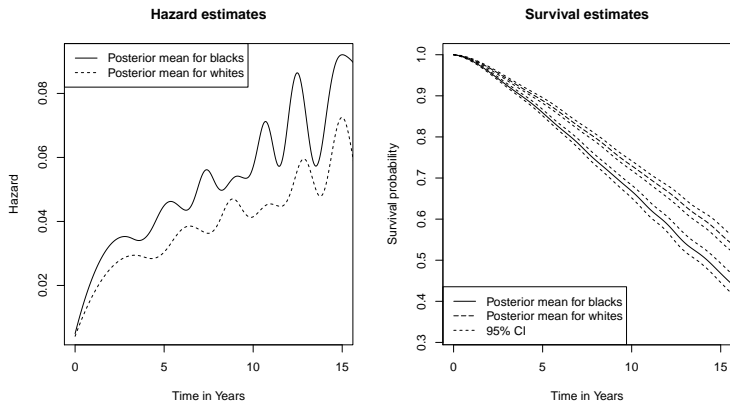


Figure: Hazard and survival for black patients (solid line) and white patients; baseline covariates.

Interpretation for race effect

- Based on reduced models, white South Carolina subjects diagnosed with prostate cancer in live 22% longer ($e^{0.20} \approx 1.22$) than black patients (95% CI is 18% to 25%), fixing age, stage, and SEER stage.
- Cox said “...*the physical or substantive basis for...proportional hazards models...is one of its weaknesses...*” and goes on to suggest that “...*accelerated failure time models are in many ways more appealing because of their quite direct physical interpretation.*”
- The SCCCR analysis showed that the main covariate of interest, race, is best modeled as an AFT effect.
- Survival probabilities for black patients are significantly lower than those for white patients when other factors are fixed at the same levels.

More interpretation

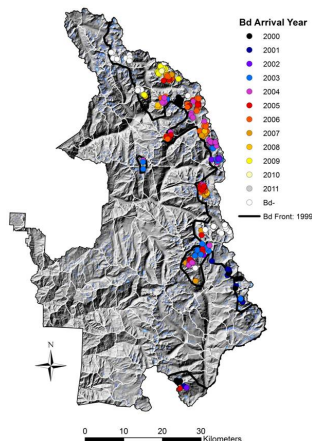
- Decreasing age by one year increases survival time by 5.4%.
- Hazard of dying increases 46% for poorly or undifferentiated grades vs. well or moderately differentiated, holding age, race, and SEER stage constant.
- SEER stage has general ExtH effects, $e^{2.77} \approx 16$ (AH) and $e^{1.21} \approx 3.4$ (PH). Those with distant stage are at least three times worse in one-sixteenth of the time as those with localized or regional.
- Marital status essentially has PH interpretation; single (including widowed or separated) subjects are $e^{0.33} \approx 1.39$ times more likely to die at any instant than married.

Extinction of mountain yellow-legged frog

- Frogs and other amphibians have been dying off in large numbers since the 1980s because of a deadly fungus called *Batrachochytrium dendrobatidis*, or Bd.
- Dr. Knapp has been studying the amphibian declines for the past decade at Sierra Nevada Aquatic Research Laboratory; he has hiked thousands of miles and surveyed hundreds of frog populations in Sequoia-Kings Canyon National Park collecting the data by hand.
- As with the SCCCR data, proportional hazards grossly violated.
- Instead of semiparametric, pursue nonparametric $F_{\mathbf{x}_i}$; not able to use frailties.

The Frog Data (2000-2011)

- Contains 309 frog populations. Each was followed up until infection or being censored (10% censoring).
- Response T_i is time to Bd infection. (i.e. Bd arrival year – baseline year).
- Main covariates:
 - $x_{i1} \in \{0, 1\}$ is whether or not Bd has been found in the watershed.
 - x_{i2} is straight-line distance to the nearest Bd location.
- Populations near each other tend to become infected at about the same time.



LDDPM model and Spatial Extension

- LDDPM (De Iorio et al., 2009; Jara et al., 2010): $Z_i = \log T_i$ given \mathbf{x}_i follows mixture model

$$F_{\mathbf{x}_i}(z) = \int \Phi\left(\frac{z - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) dG(\boldsymbol{\beta}, \sigma^2),$$

where G follows Dirichlet Process (DP) prior:

$$G \sim DP(\alpha, G_0).$$

- Countable mixture of parametric linear models

$$F_{\mathbf{x}_i} = \sum_{j=1}^{\infty} w_j N(\mathbf{x}_i' \boldsymbol{\beta}_j, \sigma_j^2).$$

- As before, take $Y_i = \Phi^{-1}\{F_{\mathbf{x}_i}(\log T_i)\}$ and $\mathbf{Y} \sim N_n(\mathbf{0}, \boldsymbol{\Gamma})$.
- $\boldsymbol{\Gamma}_\theta$ used for capturing spatial dependence;
 $\gamma_{ij} = \theta_1 \exp\{-\theta_2 \|\mathbf{s}_i - \mathbf{s}_j\|\} + (1 - \theta_1) I\{\mathbf{s}_i = \mathbf{s}_j\}.$

MCMC Overview

- Truncated stick-breaking representation
 $G = \sum_{i=1}^N [v_i \prod_{j<i} (1 - v_j)] \delta_{\beta_j, \sigma_j^2}$ where
 $v_1, \dots, v_{N-1} \stackrel{iid}{\sim} \text{beta}(1, \alpha)$, $v_N = 1$, and $(\beta_j, \sigma_j^2) \stackrel{iid}{\sim} G_0$.
- G parameters updated based on a M-H proposal from blocked Gibbs sampler (Ishwaran and James, 2001).
- The latent censored t_i updated via M-H sampler.
- Delayed rejection (Tierney and Mira, 1999) used for several parameters; helps sampler not get “stuck.”
- Correlation parameters θ are updated using adaptive M-H (Haario et al., 2001).
- For large n , the inversion of the $n \times n$ matrix \mathbf{C} substantially sped up using a full scale approximation (FSA) (Sang and Huang, 2012).

Frog Data: Inference on Spatial Correlation

- Posterior mean $\hat{\theta}_1 = 0.9937$.
- Posterior mean $\hat{\theta}_2 = 0.0866$, indicating the correlation decays by $1 - \exp\{-0.0866(1)\} = 8\%$ for every 1-km increase in distance and $1 - \exp\{-0.0866(10)\} = 58\%$ for every 10-km increase in distance.

Table: Posterior summary statistics for the spatial correlation parameters

| Par. | Mean | Median | Std. dev. | 95% HPD Interval |
|------------|--------|--------|-----------|------------------|
| θ_1 | 0.9937 | 0.9941 | 0.0029 | (0.9879, 0.9988) |
| θ_2 | 0.0866 | 0.0841 | 0.0211 | (0.0493, 0.1297) |

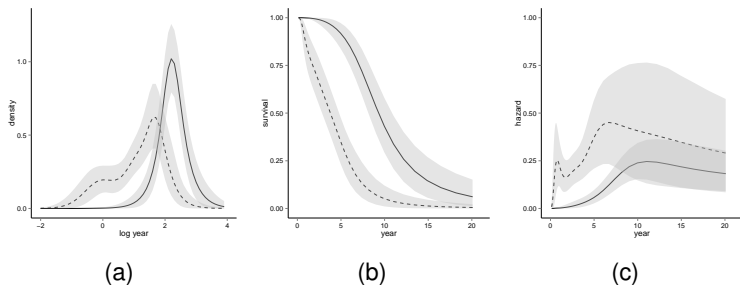


Figure: Fitted marginal densities, survival curves, and hazard curves w/ 90% CI for high versus low value of bddist when bdwater is equal to 0; bddist=95% and bddist=5% quantiles are solid and dashed lines.

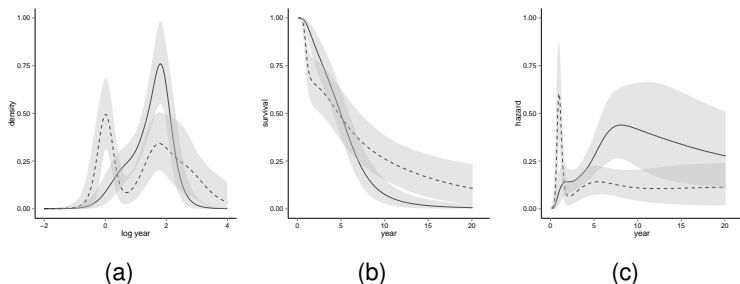


Figure: Fitted marginal densities, survival curves, and hazard curves w/ 90% CI for $bdwater=0$ versus $bdwater=1$ when $bddist$ is equal to population mean of 2.7 km; results for $bdwater=0$ and $bdwater=1$ are solid and dashed lines.

Frog Data: Spatial Prediction

Spatial map for the transformed process

$$z(\mathbf{s}) = \Phi^{-1} \{ F_{\mathbf{x}(\mathbf{s})}(\log T(\mathbf{s})|G) \}.$$

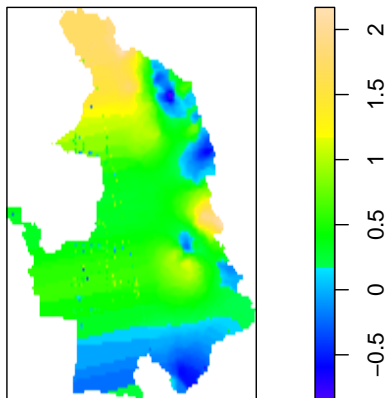


Figure: Predictive spatial map across \mathcal{D} .

Which is better, copula or frailty?

| LPML | Model |
|------|-------------------|
| -276 | LDDPM-copula |
| -304 | PH-copula |
| -632 | LDDPM-independent |
| -705 | PH-independent |
| -703 | PH-frailty |

LDDPM copula model better than PH copula model. However, PH copula better than LDDPM without copula. Modeling via copula grossly improves predictive performance of the models. Frailty improves PH model only slightly.

Remarks

- Proposed Bayesian spatial copula approaches to estimate survival curves semiparametrically (ExtH model) and nonparametrically (LDDPM) while allowing for spatial dependence, leading to high predictive accuracy.
- Implementation of simpler semiparametric models such as proportional odds focus of current research, both frailty and copula.
- Thanks to my co-authors Li Li, Haiming Zhou, Roland Knapp, and Jiajia Zhang. Thanks for the invitation!
- Papers based on this work are available; email if interested.