## GPSRL: Learning Semi-Parametric Bayesian Survival Rule Lists from Heterogeneous Patient Data

#### Ameer Hamza Shakur<sup>1</sup>, Xiaoning Qian<sup>2</sup>, Zhangyang Wang<sup>3</sup>, Bobak Mortazavi<sup>2</sup>, Shuai Huang<sup>1</sup>

<sup>1</sup>University of Washington, Seattle, WA <sup>2</sup>Texas A&M University, College Station, TX <sup>3</sup>The University of Texas at Austin, Austin, TX

#### Motivation

Patient health data is often collected from a *heterogeneous* population of patients

#### Motivation

- Patient health data is often collected from a *heterogeneous* population of patients
- Standard survival models focus on average effect of the covariates on survival outcomes

#### Motivation

- Patient health data is often collected from a *heterogeneous* population of patients
- Standard survival models focus on average effect of the covariates on survival outcomes
- Advances in sensing technology have provided opportunities to further model *heterogeneity* as well as *non-linearity* of the survival risk.

# **Survival Analysis**

• Cox proportional hazards model:

$$h(t \mid \mathbf{x}_i) = h_0(t) \exp(g(\mathbf{x}_i))$$
$$g(\mathbf{x}_i) = \mathbf{w}' \mathbf{x}_i$$

• Accelerated failure time (AFT) model:

 $T_{i} = \exp(g(\mathbf{x}_{i}))T_{0}$ log  $T_{i} = g(\mathbf{x}_{i}) + \log(T_{0}) = g(\mathbf{x}_{i}) + \epsilon$ assigning a distribution to  $\epsilon$ ,  $\epsilon \sim \operatorname{logistic}(\beta^{-1}) \implies T_{i} \sim \operatorname{log-logistic}(\alpha_{i}, \beta)$ where,  $\alpha_{i} = \exp(g(\mathbf{x}_{i})) \quad g(\mathbf{x}_{i}) = \mathbf{w}'\mathbf{x}_{i}$ 

# **Survival Analysis**

• Cox proportional hazards model:

• Accelerated failure time (AFT) model:  

$$T_{i} = \exp(g(\mathbf{x}_{i}))T_{0}$$

$$\log T_{i} = g(\mathbf{x}_{i}) + \log(T_{0}) = g(\mathbf{x}_{i}) + \epsilon$$
assigning a distribution to  $\epsilon$ ,  
 $\epsilon \sim \operatorname{logistic}(\beta^{-1}) \implies T_{i} \sim \operatorname{log-logistic}(\alpha_{i}, \beta)$ 
where,  $\alpha_{i} = \exp(g(\mathbf{x}_{i}))$ 

$$g(\mathbf{x}_{i}) = \mathbf{w}'\mathbf{x}_{i}$$

- Nonlinear effects?
- Replace linear effects g(x) with a GP f

- Nonlinear effects?
- Replace linear effects g(x) with a GP f
- GP-Cox
- GP-AFT survival model

 $\boldsymbol{\alpha} = \exp(\mathbf{f})$  where  $\mathbf{f} \sim GP(0, \mathbf{K})$ 

- Nonlinear effects?
- Replace linear effects *g(x)* with a GP **f**
- GP-Cox
- GP-AFT survival model

$$\boldsymbol{\alpha} = \exp(\mathbf{f})$$
 where  $\mathbf{f} \sim GP(0, \mathbf{K})$   $\log(\beta) \sim U(0, s)$ 

- Nonlinear effects?
- Replace linear effects g(x) with a GP f
- GP-Cox
- GP-AFT survival model

$$\boldsymbol{\alpha} = \exp(\mathbf{f})$$
 where  $\mathbf{f} \sim GP(0, \mathbf{K})$   $\log(\beta) \sim U(0, s)$ 

- but, like AFT,  $\beta$  is assumed to be constant and does not depend on co-variates
- restrictive assumption on *shape parameter* in heterogenous datasets

# Heterogeneity in GP's

- Data partitioning approaches
  - Ex. Bayesian treed partitioning, Voronoi tessellations etc.
- Trees or more complex structures used for partitioning & modelling heterogeneity



# Heterogeneity in GP's

- Data partitioning approaches
  - Ex. Bayesian treed partitioning, Voronoi tessellations etc.
- Trees or more complex structures used for partitioning & modelling heterogeneity
- Fully probabilistic and computationally demanding
- Do not address nonlinear survival models such as GP-AFT.



• Pre-mine: *R*, *a* rule set containing *K* rules

- Pre-mine: *R*, *a* rule set containing *K* rules
- Construct:  $d = \{r_1, r_2 \cdots, r_m\} \subset R$  an <u>ordered rule list</u> of size m

```
d \begin{array}{c} \textit{if } r_1 \textit{ then } t_1 \sim LL(\boldsymbol{\alpha_1}, \beta_1) \\ \textit{else if } r_2 \textit{ then } t_2 \sim LL(\boldsymbol{\alpha_2}, \beta_2) \\ & \ddots \\ \textit{else if } r_m \textit{ then } t_m \sim LL(\boldsymbol{\alpha_m}, \beta_m) \\ & \textit{else t}_0 \sim LL(\boldsymbol{\alpha_0}, \beta_0) \end{array}
```

- Pre-mine: *R*, *a* rule set containing *K* rules
- Construct:  $d = \{r_1, r_2 \cdots, r_m\} \subset R$  an <u>ordered rule list</u> of size m

```
d \begin{array}{c} \textit{if } r_1 \textit{ then } t_1 \sim LL(\boldsymbol{\alpha_1}, \beta_1) \\ \textit{else if } r_2 \textit{ then } t_2 \sim LL(\boldsymbol{\alpha_2}, \beta_2) \\ & \ddots \\ \textit{else if } r_m \textit{ then } t_m \sim LL(\boldsymbol{\alpha_m}, \beta_m) \\ & \textit{else t}_0 \sim LL(\boldsymbol{\alpha_0}, \beta_0) \end{array}
```

- Pre-mine: *R*, *a* rule set containing *K* rules
- Construct:  $d = \{r_1, r_2 \cdots, r_m\} \subset R$  an <u>ordered rule list</u> of size m

```
d \begin{array}{c} \textit{if } r_1 \textit{ then } t_1 \sim LL(\boldsymbol{\alpha_1}, \beta_1) \\ \textit{else if } r_2 \textit{ then } t_2 \sim LL(\boldsymbol{\alpha_2}, \beta_2) \\ & \cdot \\ \textit{else if } r_m \textit{ then } t_m \sim LL(\boldsymbol{\alpha_m}, \beta_m) \\ & \textit{else if } r_0 \sim LL(\boldsymbol{\alpha_0}, \beta_0) \end{array}
```

• Latent GP to model non-linearities

- Pre-mine: R, a rule set containing K rules
- Construct:  $d = \{r_1, r_2 \cdots, r_m\} \subset R$  an <u>ordered rule list</u> of size m

```
d \begin{array}{c} \textit{if } r_1 \textit{ then } t_1 \sim LL(\boldsymbol{\alpha_1}, \beta_1) \\ \textit{else if } r_2 \textit{ then } t_2 \sim LL(\boldsymbol{\alpha_2}, \beta_2) \\ & \ddots \\ \textit{else if } r_m \textit{ then } t_m \sim LL(\boldsymbol{\alpha_m}, \beta_m) \\ & \textit{else t}_0 \sim LL(\boldsymbol{\alpha_0}, \beta_0) \end{array}
```

- Latent GP to model non-linearities
- Ordered rule lists to model heterogeneity

- Pre-mine: R, a rule set containing K rules
- Construct:  $d = \{r_1, r_2 \cdots, r_m\} \subset R$  an <u>ordered rule list</u> of size m

```
d \begin{array}{c} \textit{if } r_1 \textit{ then } t_1 \sim LL(\boldsymbol{\alpha_1}, \beta_1) \\ \textit{else if } r_2 \textit{ then } t_2 \sim LL(\boldsymbol{\alpha_2}, \beta_2) \\ & \ddots \\ \textit{else if } r_m \textit{ then } t_m \sim LL(\boldsymbol{\alpha_m}, \beta_m) \\ & \textit{else t}_0 \sim LL(\boldsymbol{\alpha_0}, \beta_0) \end{array}
```

- Latent GP to model non-linearities
- Ordered rule lists to model heterogeneity
- Relaxes assumption on *shape parameter*,  $\beta$

- Pre-mine: *R*, *a* rule set containing *K* rules
- Construct:  $d = \{r_1, r_2 \cdots, r_m\} \subset R$  an <u>ordered rule list</u> of size m

```
d \begin{array}{c} \text{if } r_1 \text{ then } t_1 \sim LL(\boldsymbol{\alpha_1}, \beta_1) \\ \text{else if } r_2 \text{ then } t_2 \sim LL(\boldsymbol{\alpha_2}, \beta_2) \\ & \ddots \\ \text{else if } r_m \text{ then } t_m \sim LL(\boldsymbol{\alpha_m}, \beta_m) \\ & \text{else } t_0 \sim LL(\boldsymbol{\alpha_0}, \beta_0) \end{array}
```

- Latent GP to model non-linearities
- Ordered rule lists to model heterogeneity
- Relaxes assumption on shape parameter,  $\beta$
- Rule lists strike a balance between greedy and optimal data partitioning

• Posterior probability density

 $p(d \mid \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{X}, d)p(d).$ 

• Prior probability

$$p(d) = p(m \mid \lambda) \prod_{j=1}^{m} p(c_j \mid c_1 \cdots c_{j-1}, \eta) p(r_j \mid r_1, \cdots r_{j-1}, c_j)$$

$$p(m \mid \mathcal{R}, \lambda) = \frac{(\lambda^m / m!)}{\sum_{j=1}^{K} (\lambda^j / j!)}, \quad m = 0, 1, \dots K$$

$$p(c_j \mid c_1 \cdots c_{j-1}, \eta) = \frac{\eta^{c_j} / c_j!}{\sum_{c \in C_j} \eta^c / c!}$$

$$p(r_j \mid r_1, \cdots r_{j-1}, c_j) = \frac{1}{|\{r_i \mid r_i \in \mathcal{R}_j, c_i = c_j\}|}$$

• Log-logistic likelihood

$$p(y \mid \boldsymbol{\alpha}, \beta) = \prod_{i=1}^{M:\boldsymbol{\delta}=1} \frac{(\frac{\boldsymbol{\alpha}_i}{\beta})(\frac{y_i}{\boldsymbol{\alpha}_i})^{\beta-1}}{(1+\frac{y_i}{\boldsymbol{\alpha}_i}\beta)^2} \prod_{j=1}^{N:\boldsymbol{\delta}_j=0} \frac{1}{1+(\frac{y_i}{\boldsymbol{\alpha}_i})^{\beta}}.$$

- Marginal likelihood
  - Likelihood is not conjugate-Gaussian
  - Laplace approximation

$$\begin{aligned} q(\mathbf{f}) &\approx p(\mathbf{f} \mid \mathbf{y}, \mathbf{X}) \\ p(\mathbf{y} \mid \mathbf{X}) &\approx p_L(\mathbf{y} \mid \mathbf{X}) = \int p(\mathbf{y} \mid \mathbf{f}) q(\mathbf{f}) d\mathbf{f} \end{aligned}$$

• Log-logistic likelihood

$$p(y \mid \boldsymbol{\alpha}, \beta) = \prod_{i=1}^{M:\boldsymbol{\delta}=1} \frac{(\frac{\boldsymbol{\alpha}_i}{\beta})(\frac{y_i}{\boldsymbol{\alpha}_i})^{\beta-1}}{(1+\frac{y_i}{\boldsymbol{\alpha}_i}\beta)^2} \prod_{j=1}^{N:\boldsymbol{\delta}_j=0} \frac{1}{1+(\frac{y_i}{\boldsymbol{\alpha}_i})^{\beta}}.$$

- Marginal likelihood
  - Likelihood is not conjugate-Gaussian
  - Laplace approximation

$$q(\mathbf{f}) \approx p(\mathbf{f} \mid \mathbf{y}, \mathbf{X})$$

$$p(\mathbf{y} \mid \mathbf{X}) \approx p_L(\mathbf{y} \mid \mathbf{X}) = \int p(\mathbf{y} \mid \mathbf{f}) q(\mathbf{f}) d\mathbf{f}$$

$$p_L(\mathbf{y} \mid \mathbf{X}, d) = \prod_{j=0}^{m} p_L(\mathbf{y}_j \mid \mathbf{X}_j)$$
Approx. marginal likelihood

• Initial random list:  $d^0 \sim p(d)$ 

- Initial random list:  $d^0 \sim p(d)$
- At step *t* in the sequence:  $d^{t+1} \sim Q(d^t)$

• Proposal distribution: 
$$Q(d^{t+1} | d^t, \mathcal{R}) = \begin{cases} \frac{1}{(K-m^t)(m^t+1)} & \text{if a rule is added} \\ \frac{1}{m^t} & \text{if a rule is removed} \\ \frac{1}{m^t(m^t-1)} & \text{if a rule is moved} \end{cases}$$

- Initial random list:  $d^0 \sim p(d)$
- At step *t* in the sequence:  $d^{t+1} \sim Q(d^t)$

• Proposal distribution: 
$$Q(d^{t+1} \mid d^t, \mathcal{R}) = \begin{cases} \frac{1}{(K-m^t)(m^t+1)} & \text{if a rule is added} \\ \frac{1}{m^t} & \text{if a rule is removed} \\ \frac{1}{m^t(m^t-1)} & \text{if a rule is moved.} \end{cases}$$
  
• Acceptance probability:  $\pi(d^{t+1} \mid d^t) = \min\left\{\frac{Q(d^{t+1}, d^t)}{Q(d^t, d^{t+1})} \frac{p_L(\mathbf{y} \mid \mathbf{X}, d^{t+1})p(d^{t+1})}{p_L(\mathbf{y} \mid \mathbf{X}, d^t)p(d^t)}, 1\right\}.$ 

- Initial random list:  $d^0 \sim p(d)$
- At step *t* in the sequence:  $d^{t+1} \sim Q(d^t)$

• Proposal distribution: 
$$Q(d^{t+1} \mid d^t, \mathcal{R}) = \begin{cases} \frac{1}{(K-m^t)(m^t+1)} & \text{if a rule is added} \\ \frac{1}{m^t} & \text{if a rule is removed} \\ \frac{1}{m^t(m^t-1)} & \text{if a rule is moved.} \end{cases}$$

- Acceptance probability:  $\pi(d^{t+1} \mid d^t) = \min\left\{\frac{Q(d^{t+1}, d^t)}{Q(d^t, d^{t+1})} \frac{p_L(\mathbf{y} \mid \mathbf{X}, d^{t+1})p(d^{t+1})}{p_L(\mathbf{y} \mid \mathbf{X}, d^t)p(d^t)}, 1\right\}.$
- MCMC sequence  $d^0, d^1 \cdots d^n \rightarrow p(d \mid \mathbf{X}, \mathbf{y})$

# **Experiments – Performance evaluation**

- Two datasets
  - Synthetic
  - MIMIC-III Sepsis

# **Experiments – Performance evaluation**

- Two datasets
  - Synthetic
  - MIMIC-III Sepsis
- Performance evaluation:
  - Concordance Index (C-Index)
  - Negative log probability density (NLPD)

$$\text{NLPD}(\mathbf{y}^*, \mathbf{X}^*) = \frac{\sum_{i=1}^{N} p(y_i^* | \mathbf{x}_i^*, \mathbf{X}, \mathbf{y})}{N}$$

# **Experiments – Performance evaluation**

- Two datasets
  - Synthetic
  - MIMIC-III Sepsis
- Performance evaluation:
  - Concordance Index (C-Index)
  - Negative log probability density (NLPD)

$$\text{NLPD}(\mathbf{y}^*, \mathbf{X}^*) = \frac{\sum_{i=1}^{N} p(y_i^* | \mathbf{x}_i^*, \mathbf{X}, \mathbf{y})}{N}$$

#### Results – Synthetic Data

Table 1: Estimate of ordered rule list d from the posterior

Rules	
$r_1$	$x_3 \le 0.259$
$r_2$	$x_4 > 0.596 \& x_3 > 0.196$
$r_3$	$x_4 \le 0.677 \ \& \ x_4 > 0.192$

#### **Results – Synthetic Data**

Table 1: Estimate of ordered rule list d from the posterior

Rules	
$r_1$	$x_3 \le 0.259$
$r_2$	$x_4 > 0.596 \& x_3 > 0.196$
$r_3$	$x_4 \le 0.677 \& x_4 > 0.192$



#### Results – Synthetic Data

Table 1: Estimate of ordered rule list d from the posterior

Rules	
$r_1$	$x_3 \le 0.259$
$r_2$	$x_4 > 0.596 \& x_3 > 0.196$
$r_3$	$x_4 \le 0.677 \ \& \ x_4 > 0.192$



# Experiments – MIMIC III

• MIMIC-III, a database comprising anonymized information relating to patients admitted to BIDMC<sup>2</sup> during 2001-02 and diagnosed with sepsis.

Table 1: Estimate of ordered rule list $d$ from the posterior		
Rules		
$r_1$	$\operatorname{artpH-(mean)} <= 7.249$	
$r_2$	$O2sat-(sd) \le 4.66 \& diaBP-(mean) \le 61.26$	
$r_3$	$O2sat-(sd) \le 4.8$	

#### Experiments – MIMIC III

• MIMIC-III, a database comprising anonymized information relating to patients admitted to BIDMC<sup>2</sup> during 2001-02 and diagnosed with sepsis.

Table 1: Estimate of ordered rule list $d$ from the posterior		
Rules		
$r_1$	$\operatorname{artpH-(mean)} <= 7.249$	
$r_2$	$O2sat-(sd) \le 4.66 \& diaBP-(mean) \le 61.26$	
$r_3$	$O2sat-(sd) \le 4.8$	



2. Beth Israel Beth Israel Deaconess Medical Center (BIDMC), Boston, MA

#### Experiments – MIMIC III

• MIMIC-III, a database comprising anonymized information relating to patients admitted to BIDMC<sup>2</sup> during 2001-02 and diagnosed with sepsis.

Table 1: Estimate of ordered rule list $d$ from the posterior		
Rules		
$r_1$	$\operatorname{artpH-(mean)} <= 7.249$	
$r_2$	$O2sat-(sd) \le 4.66 \& diaBP-(mean) \le 61.26$	
$r_3$	$O2sat-(sd) \le 4.8$	



2. Beth Israel Beth Israel Deaconess Medical Center (BIDMC), Boston, MA





 Gaussian Process Survival Rule Lists (GPSRL) to model heterogeneity in survival data-sets

- Gaussian Process Survival Rule Lists (GPSRL) to model heterogeneity in survival data-sets
- Semi-parametric Bayesian framework to partition the data into subsets with different survival characteristics.

- Gaussian Process Survival Rule Lists (GPSRL) to model heterogeneity in survival data-sets
- Semi-parametric Bayesian framework to partition the data into subsets with different survival characteristics.
- Addresses some limitations of standard survival Gaussian process models

- Gaussian Process Survival Rule Lists (GPSRL) to model heterogeneity in survival data-sets
- Semi-parametric Bayesian framework to partition the data into subsets with different survival characteristics.
- Addresses some limitations of standard survival Gaussian process models
- Interpretability in the form of rules

- Gaussian Process Survival Rule Lists (GPSRL) to model heterogeneity in survival data-sets
- Semi-parametric Bayesian framework to partition the data into subsets with different survival characteristics.
- Addresses some limitations of standard survival Gaussian process models
- Interpretability in the form of rules
- Performance evaluations demonstrate the effectiveness of our model

# Thank You!

#### Synthetic data simulation

• Simulated a heterogeneous survival dataset of size N = 1000, P = 4

$$D = (\mathbf{t}, \boldsymbol{\delta}, \mathbf{X})$$
  $(\mathbf{x}_i \sim U(0, 1) \; \forall i \in 1 : N)$ 

• Event times sampled from log-logistic (LL) distribution

$$t_{i} \sim LL(\alpha_{i}, \beta_{i})$$

$$\alpha(\mathbf{x}) = I_{1}\alpha_{1}(\mathbf{x}) + I_{2}\alpha_{2}(\mathbf{x}) + I_{3}\alpha_{3}(\mathbf{x})$$

$$\beta(\mathbf{x}) = I_{1}\beta_{1}(\mathbf{x}) + I_{2}\beta_{2}(\mathbf{x}) + I_{3}\beta_{3}(\mathbf{x})$$

$$\alpha_{i}(\mathbf{x}) = a_{1}\exp\left(a_{2}\left(\sum_{k=1}^{2}\exp(a_{3}(x[k] - a_{4})^{2})\right) + \sum_{k=3}^{4}\sin(\pi\mathbf{x}[k]^{2})\right)$$

$$\beta_{i}(\mathbf{x}) = b_{1}\exp\left(\sum_{k=1}^{2}\sin(2\pi\mathbf{x}[k]^{2}) + \sum_{k=3}^{4}\cos(2\pi\mathbf{x}[k]^{2})\right),$$