Bayesian inference using Catalytic prior distributions for Cox regression models

Weihao Li

Joint work with Dongming Huang

Department of Statistics, NUS

Dec 6, 2023

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distril

Dec 6, 2023 1 / 27

• Survival analysis, often used in medical research, epidemiology.

• Survival analysis, often used in medical research, epidemiology.



э

• • = • •

Image: Image:

• Survival analysis, often used in medical research, epidemiology.



• Large dimension & insufficient sample.

< ∃ ►

Motivation

• Cox regression model: popular way in survival analysis

Motivation

• Cox regression model: popular way in survival analysis

 $\begin{array}{l} \mbox{library(survival)} \\ \mbox{coxph(Surv(time, status)} \sim \mbox{X}) \end{array}$

Motivation

• Cox regression model: popular way in survival analysis

 $\begin{array}{l} \mbox{library(survival)} \\ \mbox{coxph(Surv(time, status)} \sim \mbox{X}) \end{array}$

• [Zhang et al., 2022] p = 200, n = 400



p/n= 0.5

Introduction of Cox regression model

Observed data: D = {(X_i, Y_i, δ_i)}ⁿ_{i=1}, X_i ∈ ℝ^p is the feature vector, δ_i is the censoring indicator, and Y_i = min(T_i, C_i), where T_i is the survival time for an uncensored subject and C_i is the censoring time for a censored subject.

Introduction of Cox regression model

- Observed data: D = {(X_i, Y_i, δ_i)}ⁿ_{i=1}, X_i ∈ ℝ^p is the feature vector, δ_i is the censoring indicator, and Y_i = min(T_i, C_i), where T_i is the survival time for an uncensored subject and C_i is the censoring time for a censored subject.
- Cox regression model: semi-parametric

$$L(\beta, h_0 \mid \boldsymbol{D}) = \prod_{i=1}^{n} \left\{ \exp\left(\boldsymbol{X}'_{i}\beta\right) h_0(Y_i) \right\}^{\delta_i} \exp\left\{-\exp\left(\boldsymbol{X}'_{i}\beta\right) H_0(Y_i) \right\}$$

• $h_0(t)$: baseline hazard function, unknown nuisance parameter

Introduction of Cox regression model

- Observed data: D = {(X_i, Y_i, δ_i)}ⁿ_{i=1}, X_i ∈ ℝ^p is the feature vector, δ_i is the censoring indicator, and Y_i = min(T_i, C_i), where T_i is the survival time for an uncensored subject and C_i is the censoring time for a censored subject.
- Cox regression model: semi-parametric

$$L(\beta, h_0 \mid \boldsymbol{D}) = \prod_{i=1}^{n} \left\{ \exp\left(\boldsymbol{X}'_i\beta\right) h_0(Y_i) \right\}^{\delta_i} \exp\left\{-\exp\left(\boldsymbol{X}'_i\beta\right) H_0(Y_i) \right\}$$

h₀(t) : baseline hazard function, unknown nuisance parameter
 Partial likelihood [Cox, 1975]: MPLE maximizes

$$\mathsf{PL}(\boldsymbol{\beta}) = \prod_{i=1}^{n} \left\{ \frac{\exp\left(\boldsymbol{X}_{i}^{\prime}\boldsymbol{\beta}\right)}{\sum_{i \in \mathbf{R}_{i}} \exp\left(\boldsymbol{X}_{i}^{\prime}\boldsymbol{\beta}\right)} \right\}^{\delta_{i}}$$

where $\mathbf{R}_i := \{j : Y_j \ge Y_i\}$ is the collection of index for those are risky at time Y_i .

Recap the problem of classical method

• When dimension is large or insufficient sample, MPLE suffers

- highly variable
- 2 not exist
- Iarge bias

Recap the problem of classical method

• When dimension is large or insufficient sample, MPLE suffers

- highly variable
- 2 not exist
- Iarge bias
- Similar problems for MLE

$$\hat{\boldsymbol{ heta}}_{MLE} = rg\max_{\boldsymbol{ heta}} \left(\prod_{i}^{n} f(\boldsymbol{Y}_{i} \mid \boldsymbol{X}_{i}, \boldsymbol{ heta}) \right)$$

- large variability for large p
- MLE may not exist, e.g.: Logistic regression
- large bias

- Existing strategy for MLE:
 - [Huang et al., 2020] proposed class of automatic priors called **catalytic prior** for GLM to provide <u>stable</u> estimation in high dimensional model.

- Existing strategy for MLE:
 - [Huang et al., 2020] proposed class of automatic priors called **catalytic prior** for GLM to provide <u>stable</u> estimation in high dimensional model.
 - Increase 'sample size' by generating synthetic data $\{(Y_i^*, X_i^*)\}_{i=1}^M$.

$$Obs: \{(Y_i, X_i)\}_{i=1}^n \qquad Syn: \{(Y_i^*, X_i^*)\}_{i=1}^M$$

- Existing strategy for MLE:
 - [Huang et al., 2020] proposed class of automatic priors called **catalytic prior** for GLM to provide <u>stable</u> estimation in high dimensional model.
 - Increase 'sample size' by generating synthetic data $\{(Y_i^*, X_i^*)\}_{i=1}^M$.

$$Obs: \{(Y_i, X_i)\}_{i=1}^n \qquad Syn: \{(Y_i^*, X_i^*)\}_{i=1}^M$$

• Catalytic prior based on synthetic data

$$\pi_{cat,M}(\boldsymbol{\theta} \mid \tau) \propto \left(\prod_{i=1}^{M} f(Y_{i}^{*} \mid \boldsymbol{X}_{i}^{*}, \boldsymbol{\theta})\right)^{\tau/M}$$

- Existing strategy for MLE:
 - [Huang et al., 2020] proposed class of automatic priors called **catalytic prior** for GLM to provide <u>stable</u> estimation in high dimensional model.
 - Increase 'sample size' by generating synthetic data $\{(Y_i^*, X_i^*)\}_{i=1}^M$.

$$Obs: \{(Y_i, X_i)\}_{i=1}^n \qquad Syn: \{(Y_i^*, X_i^*)\}_{i=1}^M$$

• Catalytic prior based on synthetic data

$$\pi_{cat,M}(\boldsymbol{\theta} \mid \tau) \propto \left(\prod_{i=1}^{M} f\left(\boldsymbol{Y}_{i}^{*} \mid \boldsymbol{X}_{i}^{*}, \boldsymbol{\theta}\right)\right)^{\tau/M}$$

• Posterior distribution:

$$\pi(\boldsymbol{\theta} \mid \{(Y_i, \boldsymbol{X}_i)\}_{i=1}^n) \propto \left(\prod_i^n f(Y_i \mid \boldsymbol{X}_i, \boldsymbol{\theta})\right) \pi_{cat, M}(\boldsymbol{\theta} \mid \tau)$$
$$\propto \exp\left\{\sum_{i=1}^n \log\left(f\left(Y_i \mid \boldsymbol{X}_i, \boldsymbol{\theta}\right)\right) + \frac{\tau}{M} \sum_{i=1}^M \log\left(f\left(Y_i^* \mid \boldsymbol{X}_i^*, \boldsymbol{\theta}\right)\right)\right\}$$

- Experts focus on several important covariates.
 - E.g.: thyroid-related diseases-gender, emotional state, age, etc.

¹Other type of simpler models are possible

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distribution

Experts focus on several important covariates.
 E.g.: thyroid-related diseases-gender, emotional state, age, etc.
 Only fit simpler model with subset of covariates {X₁, X₂, X₃}¹.
 → ψ̂ = (β̂₁, β̂₂, β̂₃) → predictive distribution g_{*} (y | x)

¹Other type of simpler models are possible

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distribution

- Experts focus on several important covariates.
 E.g.: thyroid-related diseases-gender, emotional state, age, etc.
 Only fit simpler model with subset of covariates {X₁, X₂, X₃}¹.
 → ψ̂ = (β̂₁, β̂₂, β̂₃) → predictive distribution g_{*} (y | x)
- Create fake data to "increase" sample size. M synthetic data points $\{(Y_i^*, X_i^*)\}_{i=1}^M$ can be generated according to following strategy:

$$\boldsymbol{X}_{i}^{*} \stackrel{i.i.d.}{\sim} Q(\boldsymbol{x}), \quad Y_{i}^{*} \mid \boldsymbol{X}_{i}^{*} \sim g_{*}\left(y \mid \boldsymbol{X}_{i}^{*}\right)$$

¹Other type of simpler models are possible

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distri

- Experts focus on several important covariates.
 E.g.: thyroid-related diseases-gender, emotional state, age, etc.
 Only fit simpler model with subset of covariates {X₁, X₂, X₃}¹.
 → ψ̂ = (β̂₁, β̂₂, β̂₃) → predictive distribution g_{*} (y | x)
- Create fake data to "increase" sample size. M synthetic data points $\{(Y_i^*, X_i^*)\}_{i=1}^M$ can be generated according to following strategy:

$$\boldsymbol{X}_{i}^{*} \stackrel{i.i.d.}{\sim} Q(\boldsymbol{x}), \quad Y_{i}^{*} \mid \boldsymbol{X}_{i}^{*} \sim g_{*}\left(y \mid \boldsymbol{X}_{i}^{*}\right)$$

• Now we are able to define catalytic prior for ${m eta}$

$$\pi_{\textit{cat},\textit{M}}(oldsymbol{ heta} \mid au) \propto \left(\prod_{i=1}^{\textit{M}} f\left(Y_{i}^{*} \mid oldsymbol{X}_{i}^{*}, oldsymbol{ heta}
ight)
ight)^{ au/\textit{M}}$$

¹Other type of simpler models are possible Weihao Li (Department of Statistics, NUS) Payelian inforence using, Catalytic proc

Generation of synthetic survival data

- Synthetic covariates: $X_i^* \stackrel{i.i.d.}{\sim} Q(x)$
- Synthetic survival time:
 - Simpler predictive model g(y | X, δ, ψ) with ψ can be stably fitted from the observed data {(X_i, Y_i, δ_i)}ⁿ_{i=1}, denote the predictive distribution as g_{*}(y | x, {(X_i, Y_i, δ_i)}ⁿ_{i=1})
 - **2** $Y_i^* | \mathbf{X}_i^* \sim g_* (y | \mathbf{X}_i^*, \{(\mathbf{X}_i, Y_i, \delta_i)\}_{i=1}^{n})$

Generation of synthetic survival data

- Synthetic covariates: $X_i^* \stackrel{i.i.d.}{\sim} Q(x)$
- Synthetic survival time:
 - Simpler predictive model g(y | X, δ, ψ) with ψ can be stably fitted from the observed data {(X_i, Y_i, δ_i)}ⁿ_{i=1}, denote the predictive distribution as g_{*}(y | x, {(X_i, Y_i, δ_i)}ⁿ_{i=1})
 - **2** $Y_i^* \mid X_i^* \sim g_* (y \mid X_i^*, \{(X_i, Y_i, \delta_i)\}_{i=1}^{n})$

Remark: all synthetic survival subjects are uncensored.

Generation of synthetic survival data

- Synthetic covariates: $\boldsymbol{X}_{i}^{*} \stackrel{i.i.d.}{\sim} Q(\boldsymbol{x})$
- Synthetic survival time:
 - Simpler predictive model g(y | X, δ, ψ) with ψ can be stably fitted from the observed data {(X_i, Y_i, δ_i)}ⁿ_{i=1}, denote the predictive distribution as g_{*}(y | x, {(X_i, Y_i, δ_i)}ⁿ_{i=1})

2
$$Y_i^* \mid \mathbf{X}_i^* \sim g_* \left(y \mid \mathbf{X}_i^*, \{ (\mathbf{X}_i, Y_i, \delta_i) \}_{i=1}^{n} \right)$$

Remark: all synthetic survival subjects are uncensored.

Example:

Simpler model:

$$L\left(\psi \mid \{(\boldsymbol{X}_{i}, Y_{i}, \delta_{i})\}_{i=1}^{n}\right) = \prod_{i=1}^{n} \psi^{\delta_{i}} \exp\left\{-\psi Y_{i}\right\}, \quad \psi > 0.$$

2 $Y_i^* \mid \boldsymbol{X}_i^* \sim \operatorname{Exp}(\hat{\psi})$

$$\boldsymbol{D} = \{(\boldsymbol{X}_{i}, Y_{i}, \delta_{i})\}_{i=1}^{n}, \boldsymbol{D}^{*} = \{(\boldsymbol{X}_{i}^{*}, Y_{i}^{*})\}_{i=1}^{M}$$

$$\pi_{\textit{cat,both}}(oldsymbol{eta},h_0) \propto L(oldsymbol{eta},h_0(\cdot) \mid oldsymbol{D}^*)^{rac{ au}{M}}$$

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distri

$$\begin{split} \boldsymbol{D} &= \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\\ &\pi_{cat, both}(\boldsymbol{\beta}, h_0) \propto L(\boldsymbol{\beta}, h_0(\cdot) \mid \boldsymbol{D}^*)^{\frac{\tau}{M}} \end{split}$$

Posterior distribution:

 $\pi_{post,both}(m{eta},h_0) \propto L(m{eta},h_0(\cdot)\mid m{D}) \cdot L(m{eta},h_0(\cdot)\mid m{D}^*)^{ au_{\overline{M}}}$

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distri

$$\boldsymbol{D} = \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M$$

$$\pi_{\textit{cat,both}}(oldsymbol{eta},h_0) \propto L(oldsymbol{eta},h_0(\cdot) \mid oldsymbol{D}^*)^{rac{1}{M}}$$

• Posterior distribution:

$$\pi_{\textit{post,both}}(\boldsymbol{eta},h_0) \propto L(\boldsymbol{eta},h_0(\cdot)\mid \boldsymbol{D}) \cdot L(\boldsymbol{eta},h_0(\cdot)\mid \boldsymbol{D}^*)^{\frac{ au}{M}}$$

Conceptually valid but not able to be incorporated into standard Bayesian Cox model.

Reason: β and h_0 are prior-dependent

$$\boldsymbol{D} = \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M$$

$$\pi_{\textit{cat,both}}(oldsymbol{eta},h_0) \propto L(oldsymbol{eta},h_0(\cdot) \mid oldsymbol{D}^*)^{rac{1}{M}}$$

• Posterior distribution:

$$\pi_{post,both}(\boldsymbol{\beta},h_0) \propto L(\boldsymbol{\beta},h_0(\cdot) \mid \boldsymbol{D}) \cdot L(\boldsymbol{\beta},h_0(\cdot) \mid \boldsymbol{D}^*)^{\frac{\tau}{M}}$$

Conceptually valid but not able to be incorporated into standard Bayesian Cox model.

Reason: β and h_0 are prior-dependent

Efficient point estimation for β: profile h₀(·) out
 [Murphy and Van der Vaart, 2000]

Algorithm: Compute Weighted Mixture estimator in R Require: R package survival

Input: $D = \{(X_i, Y_i, \delta_i)\}_{i=1}^n, D^* = \{(X_i^*, Y_i^*)\}_{i=1}^M, \tau > 0$

Algorithm: Compute Weighted Mixture estimator in R **Require: R** package survival **Input:** $\boldsymbol{D} = \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M, \tau > 0$

 $\begin{array}{l} \textbf{O} \quad \text{Combine data:} \\ \tilde{\mathbb{X}} \leftarrow \text{rbind} \left(\mathbb{X}, \mathbb{X}^*\right) \\ \tilde{\boldsymbol{Y}} \leftarrow c\left(\boldsymbol{Y}, \boldsymbol{Y}^*\right) \\ \tilde{\boldsymbol{\delta}} \leftarrow c(\boldsymbol{\delta}, \boldsymbol{1}) \end{array}$

Algorithm: Compute Weighted Mixture estimator in R Require: R package survival Input: $D = \{(X_i, Y_i, \delta_i)\}_{i=1}^n, D^* = \{(X_i^*, Y_i^*)\}_{i=1}^M, \tau > 0$

 $\begin{array}{l} \bullet \quad \text{Combine data:} \\ \tilde{\mathbb{X}} \leftarrow \mathsf{rbind}\left(\mathbb{X},\mathbb{X}^*\right) \\ \tilde{\boldsymbol{Y}} \leftarrow c\left(\boldsymbol{Y},\boldsymbol{Y}^*\right) \\ \tilde{\boldsymbol{\delta}} \leftarrow c(\boldsymbol{\delta},\boldsymbol{1}) \end{array}$

2 Compute weight vector: $\tilde{\mathbf{w}} \leftarrow c(\operatorname{rep}(1, n), \operatorname{rep}(\frac{\tau}{M}, M))$

Algorithm: Compute Weighted Mixture estimator in R Require: R package survival

- **Input:** $\boldsymbol{D} = \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M, \tau > 0$
 - $\begin{array}{l} \textbf{O} \quad \text{Combine data:} \\ \tilde{\mathbb{X}} \leftarrow \text{rbind} \left(\mathbb{X}, \mathbb{X}^*\right) \\ \tilde{\boldsymbol{Y}} \leftarrow c\left(\boldsymbol{Y}, \boldsymbol{Y}^*\right) \\ \tilde{\boldsymbol{\delta}} \leftarrow c(\boldsymbol{\delta}, \boldsymbol{1}) \end{array}$
 - **2** Compute weight vector: $\tilde{\mathbf{w}} \leftarrow c(\operatorname{rep}(1, n), \operatorname{rep}(\frac{\tau}{M}, M))$
 - Fit proportional hazards regression model with combined data and weight vector:
 fit ← coxph(Surv(Υ̃, δ̃) ~ X̃, weights = w̃)

Algorithm: Compute Weighted Mixture estimator in R Require: R package survival

- **Input:** $\boldsymbol{D} = \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M, \tau > 0$
 - $\begin{array}{l} \bullet \quad \text{Combine data:} \\ \tilde{\mathbb{X}} \leftarrow \mathsf{rbind}\left(\mathbb{X},\mathbb{X}^*\right) \\ \tilde{\boldsymbol{Y}} \leftarrow c\left(\boldsymbol{Y},\boldsymbol{Y}^*\right) \\ \tilde{\boldsymbol{\delta}} \leftarrow c(\boldsymbol{\delta},\boldsymbol{1}) \end{array}$
 - **2** Compute weight vector: $\tilde{\mathbf{w}} \leftarrow c(\operatorname{rep}(1, n), \operatorname{rep}(\frac{\tau}{M}, M))$
 - Fit proportional hazards regression model with combined data and weight vector:
 fit ← coxph(Surv(Υ̃, δ̃) ~ X̃, weights = w̃)

$$\hat{\boldsymbol{\beta}}_{WM,\tau} = \operatorname{coef}(\mathsf{fit})$$

Algorithm: Compute Weighted Mixture estimator in R Require: R package survival

- **Input:** $\boldsymbol{D} = \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n, \boldsymbol{D}^* = \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M, \tau > 0$
 - $\begin{array}{l} \textbf{O} \quad \text{Combine data:} \\ \tilde{\mathbb{X}} \leftarrow \mathsf{rbind}\left(\mathbb{X},\mathbb{X}^*\right) \\ \tilde{\boldsymbol{Y}} \leftarrow c\left(\boldsymbol{Y},\boldsymbol{Y}^*\right) \\ \tilde{\boldsymbol{\delta}} \leftarrow c(\boldsymbol{\delta},\boldsymbol{1}) \end{array}$
 - **2** Compute weight vector: $\tilde{\mathbf{w}} \leftarrow c(\operatorname{rep}(1, n), \operatorname{rep}(\frac{\tau}{M}, M))$
 - Fit proportional hazards regression model with combined data and weight vector:

$$\mathsf{fit} \leftarrow \mathsf{coxph}(\mathsf{Surv}(\widetilde{\boldsymbol{Y}}, \widetilde{\boldsymbol{\delta}}) \sim \widetilde{\mathbb{X}}, \mathsf{weights} = \widetilde{\boldsymbol{\mathsf{w}}})$$

(4)
$$\hat{oldsymbol{eta}}_{WM, au} = \operatorname{coef}(\operatorname{fit})$$

Property:

• $\hat{\beta}_{WM,\tau}$ is consistent when p is fixed and $\tau = o(n)$.

- Previous catalytic prior for GLM rely on known likelihood function.
- Issue: $h_0(t)$ is unknown nuisance parameter in Cox model

- Previous catalytic prior for GLM rely on known likelihood function.
- Issue: $h_0(t)$ is unknown nuisance parameter in Cox model

 $\pi_{cat,both}(\boldsymbol{\beta}, h_0) \propto L(\boldsymbol{\beta}, h_0(\cdot) \mid \boldsymbol{D}^*)^{\frac{T}{M}} \\ = \left[\prod_{i=1}^{M} \left\{ \exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) h_0(\boldsymbol{Y}_i^*) \right\} \exp\left\{ -\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) \int_0^{\boldsymbol{Y}_i^*} h_0(s) ds \right\} \right]^{\tau/M}$

- Previous catalytic prior for GLM rely on known likelihood function.
- Issue: $h_0(t)$ is unknown nuisance parameter in Cox model

$$\pi_{cox,cat}(\boldsymbol{\beta} \mid \tau) \propto \mathcal{L}\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M} \\ = \left[\prod_{i=1}^M \left\{\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) h_0^+\right\} \exp\left\{-\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) Y_i^* h_0^+\right\}\right]^{\tau/M}$$

- Previous catalytic prior for GLM rely on known likelihood function.
- Issue: $h_0(t)$ is unknown nuisance parameter in Cox model

$$\pi_{cox,cat}(\boldsymbol{\beta} \mid \tau) \propto \mathcal{L}\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M} \\ = \left[\prod_{i=1}^M \left\{\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) h_0^+\right\} \exp\left\{-\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) Y_i^* h_0^+\right\}\right]^{\tau/M}$$

- User-specific surrogate baseline hazard constant $h_0^+ > 0$
 - User-Defined constant: constant hazard, hazard with domain knowledge
 - 2 Computed in a data driven way.

- Previous catalytic prior for GLM rely on known likelihood function.
- Issue: $h_0(t)$ is unknown nuisance parameter in Cox model

$$\pi_{cox,cat}(\boldsymbol{\beta} \mid \tau) \propto \mathcal{L}\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M} \\ = \left[\prod_{i=1}^M \left\{\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) h_0^+\right\} \exp\left\{-\exp\left(\boldsymbol{X}_i^{*\prime}\boldsymbol{\beta}\right) Y_i^* h_0^+\right\}\right]^{\tau/M}$$

- User-specific surrogate baseline hazard constant $h_0^+ > 0$
 - User-Defined constant: constant hazard, hazard with domain knowledge
 - 2 Computed in a data driven way.

 h_0^+ acts merely as a surrogate for the nuisance component to facilitate the construction of our catalytic prior. It does not need to be correctly specified or unbiased.

• Example on
$$h_0^+$$
:
• $h_0^+ = \hat{\psi}$
 $\hat{\psi} = \arg \max_{\psi} \left(\prod_{i=1}^n \psi^{\delta_i} \exp \{-\psi Y_i\} \right), \quad \psi > 0.$
• $h_0^+ = 1/\bar{Y}$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Bayesian Cox model with catalytic prior

$$\pi_{cox,cat}(\beta \mid \tau) = L\left(\beta, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M}$$
(1)

Procedure:

 Standard Bayesian Cox model with L(β, h₀) as likelihood and default prior for parameter π₀(β) · π(h₀)

$$\pi_{\textit{post}}(m{eta}, h_0 \mid m{D}) \propto L(m{eta}, h_0) \cdot \pi_0(m{eta}) \cdot \pi(h_0)$$

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distril

Bayesian Cox model with catalytic prior

$$\pi_{cox,cat}(\beta \mid \tau) = L\left(\beta, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M}$$
(1)

Procedure:

 Standard Bayesian Cox model with L(β, h₀) as likelihood and default prior for parameter π₀(β) · π(h₀)

$$\pi_{\textit{post}}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) \propto L(\boldsymbol{\beta}, h_0) \cdot \pi_0(\boldsymbol{\beta}) \cdot \pi(h_0)$$

- Replace default prior $\pi_0(\beta)$ on β by catalytic prior Eq.(1).
- Run MCMC to collect posterior samples.

Bayesian Cox model with catalytic prior

$$\pi_{cox,cat}(\beta \mid \tau) = L\left(\beta, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M}$$
(1)

Procedure:

 Standard Bayesian Cox model with L(β, h₀) as likelihood and default prior for parameter π₀(β) · π(h₀)

$$\pi_{post}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) \propto L(\boldsymbol{\beta}, h_0) \cdot \pi_0(\boldsymbol{\beta}) \cdot \pi(h_0)$$

- Replace default prior $\pi_0(\beta)$ on β by catalytic prior Eq.(1).
- Run MCMC to collect posterior samples.

Property:

• $\pi_{cox,cat}(\beta \mid \tau)$ is proper under mild assumption.

Approximate posterior mode

$$\pi_{\textit{post}}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) = L(\boldsymbol{\beta}, h_0) \cdot \pi_{\textit{cox}, \textit{cat}}(\boldsymbol{\beta} \mid \tau) \cdot \pi(h_0)$$

Approximate posterior mode

(0)

$$\pi_{post}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) = L(\boldsymbol{\beta}, h_0) \cdot \pi_{cox, cat}(\boldsymbol{\beta} \mid \tau) \cdot \pi(h_0)$$

$$\arg \max_{\boldsymbol{\beta}} \pi_{margin}(\boldsymbol{\beta} \mid \boldsymbol{D}) = \arg \max_{\boldsymbol{\beta}} \int \pi_{post}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) dh_0$$

١

(1)

 $1(0 \rangle$

$$\pi_{\textit{post}}(\beta, h_0 \mid \boldsymbol{D}) = L(\beta, h_0) \cdot \pi_{\textit{cox}, \textit{cat}}(\beta \mid \tau) \cdot \pi(h_0)$$

arg max
$$\pi_{margin}(\beta \mid D) = \arg \max_{\beta} \int \pi_{post}(\beta, h_0 \mid D) dh_0$$

Logic: from [Sinha, 2003], when $\pi(h_0)$ is diffuse,

$$\pi_{margin}(\boldsymbol{\beta} \mid \boldsymbol{D}) \approx \textit{PL}(\boldsymbol{\beta}) \cdot \pi_{\textit{cox},\textit{cat}}(\boldsymbol{\beta} \mid \tau)$$

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distril

$$\pi_{\textit{post}}(\beta, h_0 \mid \boldsymbol{D}) = L(\beta, h_0) \cdot \pi_{\textit{cox}, \textit{cat}}(\beta \mid \tau) \cdot \pi(h_0)$$

$$\arg \max_{\boldsymbol{\beta}} \pi_{margin}(\boldsymbol{\beta} \mid \boldsymbol{D}) = \arg \max_{\boldsymbol{\beta}} \int \pi_{post}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) dh_0$$

Logic: from [Sinha, 2003], when $\pi(h_0)$ is diffuse,

$$\pi_{margin}(\boldsymbol{\beta} \mid \boldsymbol{D}) \approx PL(\boldsymbol{\beta}) \cdot \pi_{cox,cat}(\boldsymbol{\beta} \mid \tau)$$

Point estimation (Catalytic Regularized estimator)

$$\hat{\boldsymbol{\beta}}_{CR,\tau} = \arg\max_{\boldsymbol{\beta}} \left\{ \log PL(\boldsymbol{\beta}) + \log \pi_{cox,cat}(\boldsymbol{\beta} \mid \tau) \right\},$$
(2)

$$\pi_{\textit{post}}(\boldsymbol{\beta}, h_0 \mid \boldsymbol{D}) = L(\boldsymbol{\beta}, h_0) \cdot \pi_{\textit{cox}, \textit{cat}}(\boldsymbol{\beta} \mid \tau) \cdot \pi(h_0)$$

$$rg\max_{oldsymbol{eta}} \ \pi_{\textit{margin}}(oldsymbol{eta} \mid oldsymbol{D}) = rg\max_{oldsymbol{eta}} \ \int \pi_{\textit{post}}(oldsymbol{eta}, h_0 \mid oldsymbol{D}) dh_0$$

Logic: from [Sinha, 2003], when $\pi(h_0)$ is diffuse,

$$\pi_{margin}(\boldsymbol{\beta} \mid \boldsymbol{D}) \approx PL(\boldsymbol{\beta}) \cdot \pi_{cox,cat}(\boldsymbol{\beta} \mid \tau)$$

Point estimation (Catalytic Regularized estimator)

$$\hat{\boldsymbol{\beta}}_{CR,\tau} = \arg\max_{\boldsymbol{\beta}} \left\{ \log PL(\boldsymbol{\beta}) + \log \pi_{cox,cat}(\boldsymbol{\beta} \mid \tau) \right\},$$
(2)

Property:

• $\hat{\beta}_{CR,\tau}$ is consistent when p is fixed and $\tau = o(n)$.

Short summary

Method:

• Full Bayesian Procedure

$$\pi_{\text{cox,cat}}(\boldsymbol{\beta} \mid \tau) \quad = \quad L\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M}$$

• Catalytic Regularized estimator

$$\hat{oldsymbol{eta}}_{\textit{CR}, au} = rg\max_{oldsymbol{eta}} \{\log\textit{PL}(oldsymbol{eta}) + \log\pi_{\textit{cox},\textit{cat}}(oldsymbol{eta} \mid au)\}$$

• Weighted Mixture estimator

$$\hat{oldsymbol{eta}}_{WM, au}$$

Theory:

•
$$\pi_{cox,cat}(\beta \mid \tau)$$
 is proper
• $\hat{\beta}_{CR,\tau}$ and $\hat{\beta}_{WM,\tau}$ are consistent.

Observed data:

• Covariates:

$$\boldsymbol{X}_{i,j} \sim \begin{cases} \mathsf{Bernoulli}(0.1), & j = 1\\ \chi_1^2, & j = 2\\ \chi_4^2, & j = 3\\ \mathcal{N}(0,1), & 4 \le j \le p \end{cases}$$

Observed data:

• Covariates:

$$\boldsymbol{X}_{i,j} \sim \begin{cases} \text{Bernoulli}(0.1), & j = 1\\ \chi_1^2, & j = 2\\ \chi_4^2, & j = 3\\ N(0,1), & 4 \le j \le p \end{cases}$$

- n = 100, 20% censored subjects.
- The true regression coefficient vector is set to be $\beta_0 = (4, -4, 3, -3, \mathbf{1}_{p-4})/\sqrt{p}.$
- T_i independently from an exponential distribution with a rate parameter of $0.5 \exp(\mathbf{X}_i^\top \beta_0)$.

Synthetic data:

- M=1000
- Synthetic covariates: independent resampling from original covariates

Synthetic data:

- M=1000
- Synthetic covariates: independent resampling from original covariates Modification:
 - half of the sampled X_1^* will be replaced by i.i.d. random variables draw from Bernoulli(p = 0.5)—flattening
 - continuous covariate $(j \ge 2)$, half of sampled X_j^* will be replaced by i.i.d random variables draw from a normal distribution with median and interquartile range matching to those of the observed covariates.

Synthetic data:

- M=1000
- Synthetic covariates: independent resampling from original covariates Modification:
 - half of the sampled X_1^* will be replaced by i.i.d. random variables draw from Bernoulli(p = 0.5)—flattening
 - continuous covariate $(j \ge 2)$, half of sampled X_j^* will be replaced by i.i.d random variables draw from a normal distribution with median and interquartile range matching to those of the observed covariates.
- Synthetic survival time: $Y_i^* \sim \mathsf{Exp}(\hat{\psi})$

$$L(\psi \mid \{(\boldsymbol{X}_i, Y_i, \delta_i)\}_{i=1}^n) = \prod_{i=1}^n \psi^{\delta_i} \exp\left\{-\psi Y_i\right\}, \quad \psi > 0.$$

Numerical studies

Our methods: $h_0^+ = \hat{\psi}$

$$\pi_{\mathsf{cox},\mathsf{cat}}(\boldsymbol{\beta} \mid \tau) = L\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M}$$

- **(**) CRE: catalytic regularized estimator $\hat{oldsymbol{eta}}_{{\it CR}, au}$
- ② WME: weighted mixture estimator $\hat{oldsymbol{eta}}_{WM, au}$

OPM: posterior mean of MCMC sampler based on catalytic prior Alternative methods:

$$\hat{\boldsymbol{\beta}}_{\lambda} = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \left\{ \sum_{i=1}^{n} \delta_{i} \left(\boldsymbol{X}_{i}^{\top} \boldsymbol{\beta} - \log \sum_{j \in \mathcal{R}_{i}} \theta_{j} \right) - \lambda f(\boldsymbol{\beta}) \right\}$$

f(β) = 0 ⇒ MPLE: maximum partial likelihood estimator
 f(β) = ||β||₁ ⇒ Lasso
 f(β) = ||β||₂² ⇒ Ridge

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distri

$$\mathit{Deviance}(eta_0, \hat{oldsymbol{eta}}) = \left\{ \ell^{(\mathsf{test})}(oldsymbol{eta}_0) - \ell^{(\mathsf{test})}(\hat{oldsymbol{eta}})
ight\}$$

р	Methods	$\ \hat{oldsymbol{eta}}-oldsymbol{eta}_0\ ^2$	Predictive deviance
	MPLE	0.95(0.06)	19.73(1.32)
	CRE (CV)	0.63(0.04)	12.87(0.80)
	$CRE(\tau = p)$	0.86(0.02)	19.53(0.70)
	WME (CV)	0.51 (0.02)	12.44 (0.75)
20	CPM (CV)	0.79(0.04)	18.82(0.67)
	Ridge (CV)	0.58(0.03)	13.07(0.63)
	Lasso (CV)	0.75(0.04)	13.42(0.63)
	MPLE	4.25(0.21)	103.38(4.99)
	CRE (CV)	0.82(0.02)	24.09(0.84)
	$CRE(\tau = p)$	0.88(0.01)	25.63(0.77)
	WME (CV)	0.76 (0.02)	23.37 (0.82)
40	CPM (CV)	0.84(0.02)	24.63(0.76)
	Ridge (CV)	0.94(0.03)	24.53(0.87)
	Lasso (CV)	1.19(0.03)	24.06(0.92)

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distri

Real data: PBC dataset

- Data: p = 18, DPCA treatment + 17 covariates
- Full Bayesian analysis with $\tau = p$

Sampling distribution for Coefficient of DPCA



95% Credible interval: [-0.28,0.47]. Ineffectiveness of DPCA treatment has been reported [Locke III et al., 1996].

- Propose Cox catalytic prior on β and its corresponding approximate posterior mode.
- Derive weighted mixture estimator Using synthetic data to derive.
- Establish properness of prior and consistency of point estimation.
- Show proposed methods outperform classical MPLE and comparable with existing regularization method.

- Propose Cox catalytic prior on β and its corresponding approximate posterior mode.
- Derive weighted mixture estimator Using synthetic data to derive.
- Establish properness of prior and consistency of point estimation.
- Show proposed methods outperform classical MPLE and comparable with existing regularization method.

Future direction:

- Other statistical inference that involve partial likelihood or complex nuisance parameters.
- Other semi-parametric model.

```
http://arxiv.org/abs/2312.01411
```

Thank you!

Weihao Li (Department of Statistics, NUS) Bayesian inference using Catalytic prior distril

If a fully Bayesian perspective on τ is adopted, we can impose a joint prior on (τ, β) as follows. Given any two positive scalar hyperparameters α and γ , we define a joint catalytic prior for (τ, β) as

$$\pi_{\alpha,\gamma}(\tau,\boldsymbol{\beta}) \propto \Gamma_{\alpha,\gamma}(\tau) \cdot L\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)^{\tau/M}$$
(3)

where $\Gamma_{\alpha,\gamma}(\tau)$ is a function defined as

$$\Gamma_{\alpha,\gamma}(\tau) = \tau^{p+\alpha-1} e^{-\tau \left(\kappa+\gamma^{-1}\right)} \tag{4}$$

and

$$\kappa := \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{M} \log L\left(\boldsymbol{\beta}, h_0^+ \mid \{(\boldsymbol{X}_i^*, Y_i^*)\}_{i=1}^M\right)$$
(5)

The ridge estimator is defined as follows:

$$\hat{\boldsymbol{\beta}}_{\lambda} = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \left\{ \sum_{i=1}^{n} \delta_{i} \left(\boldsymbol{X}_{i}^{\top} \boldsymbol{\beta} - \log \sum_{j \in \mathcal{R}_{i}} \theta_{j} \right) - \lambda f(\boldsymbol{\beta}) \right\}$$
(6)

where $\theta_j = \exp(\mathbf{X}_j^{\top} \boldsymbol{\beta})$ and the penalty term $f(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|^2$, assuming that each entry of the covariates \mathbf{X}_i 's has been standardized to have zero mean and unit variance. The Lasso estimator is defined in a similar way as in (6) but with $f(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_1$. Both the ridge and the Lasso estimates can be computed using the glmnet package in **R** [Simon et al., 2011].

Reference I



```
Cox, D. R. (1975).
Partial likelihood.
Biometrika, 62(2):269–276.
```

Huang, D., Stein, N., Rubin, D. B., and Kou, S. C. (2020). Catalytic prior distributions with application to generalized linear models.

Proc. Natl. Acad. Sci. U.S.A., 117(22):12004-12010.

 Locke III, G. R., Therneau, T. M., Ludwig, J., Dickson, E. R., and Lindor, K. D. (1996).
 Time course of histological progression in primary biliary cirrhosis. *Hepatology*, 23(1):52–56.

```
    Murphy, S. A. and Van der Vaart, A. W. (2000).
    On profile likelihood.
    Journal of the American Statistical Association, 95(450):449–465.
```

Simon, N., Friedman, J., Hastie, T., and Tibshirani, R. (2011). Regularization Paths for Cox's Proportional Hazards Model via Coordinate Descent.

Journal of Statistical Software, 39(5):1–13.

Sinha, D. (2003).
 A Bayesian justification of Cox's partial likelihood.
 Biometrika, 90(3):629–641.

 Zhang, X., Zhou, H., and Ye, H. (2022).
 A modern theory for high-dimensional Cox regression models. arXiv preprint arXiv:2204.01161.