

KSP: Kolmogorov-Smirnov metric-based Post-Hoc Calibration for Survival Analysis

Jeongho Park¹, Daheen Kim¹, Cheoljun Kim²,
Hyungbin Park², Sangwook Kang¹, Gwangsu Kim^{2,*}

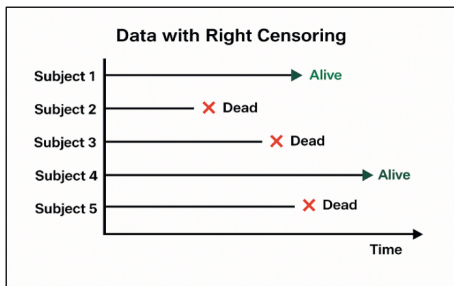
¹Department of Statistics and Data Science, Yonsei University

²Department of Statistics, Jeonbuk National University

39th Conference on Neural Information Processing Systems

Introduction

- Survival analysis aims to estimate the probability that an event (e.g., death) occurs after a given time.
- Utilizing DNN has become an essential part of the survival analysis
- A key challenge is handling censoring and balancing discrimination with calibration.



Notation

- T : event time, C : (right) censoring time
- $Y = \min \{T, C\}$: observed time
- $\delta = \mathbb{I}(T \leq C)$: censoring indicator
- \mathbf{z} : vector of covariates
- $F(\cdot | \mathbf{z})$: conditional CDF (cumulative distribution function) of T
- $S(\cdot | \mathbf{z}) = 1 - F(\cdot | \mathbf{z})$: survival function

D-calibration (Distributional calibration)

- How can we measure whether our estimated survival function is calibrated or not?
 - Use the property of CDF
 - If $T \sim F$, then $F(T) \sim \text{Unif}[0, 1]$.
 - For $T > C$, we get $F(T | \mathbf{z}) \sim \text{Unif}[F(C | \mathbf{z}), 1]$
- $\mathbb{E}_{Y, \delta, \mathbf{z}} \left[\mathbb{I}(F(Y | \mathbf{z}) \leq x) \left\{ \delta + (1 - \delta) \frac{x - F(Y | \mathbf{z})}{1 - F(Y | \mathbf{z})} \right\} \right] = x, \forall x \in [0, 1]$
- A key is how to measure the difference between both sides of the equation.

Kolmogorov-Smirnov metric

$$\mathbb{E}_{Y, \delta, \mathbf{z}} \left[\mathbb{I} (F(Y | \mathbf{z}) \leq x) \left\{ \delta + (1 - \delta) \frac{x - F(Y | \mathbf{z})}{1 - F(Y | \mathbf{z})} \right\} \right] = x$$

- Previous approaches used bin-based difference.
- We adopt the Kolmogorov-Smirnov (KS) metric.
- Let

$$\tilde{F}(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{I} (\hat{F}_{\theta}(Y_i | \mathbf{z}_i) \leq x) \left\{ \delta_i + (1 - \delta_i) \frac{x - \hat{F}_{\theta}(Y_i | \mathbf{z}_i)}{1 - \hat{F}_{\theta}(Y_i | \mathbf{z}_i)} \right\}$$

- \hat{F}_{θ} : estimated CDF
- $\text{KS-cal} = \sup_{x \in [0,1]} |\tilde{F}(x) - x|$

Convergence of KS-cal

Theorem

Under the regularity conditions,
 $\hat{F}_\theta = F$ if and only if $\sup_{x \in [0,1]} |\tilde{F}(x) - x| = o_p(1)$ as $N \rightarrow \infty$.

- Minimizing KS-cal guarantees the model is to be calibrated.
- We propose KS-cal based Post-hoc calibration (KSP).

Algorithm. KSP

-
- 1: **Input:** Estimated CDFs \hat{F}_θ , strictly monotone increasing link function $G : [0, 1] \rightarrow (-\infty, \infty)$
 - 2: Initialize parameters $a (> 0)$, b , $\alpha (> 0)$
 - 3: Sort \hat{F}_θ for computational efficiency
 - 4: **while** KS-cal not improved **do**
 - 5: Compute transformed CDF: $\hat{F}_\theta^* = \left\{ G^{-1}(a \cdot G(\hat{F}_\theta) + b) \right\}^\alpha$
 - 6: Compute KS-cal on validation set: $\max_{1 \leq j \leq N} D_j^*$, where D_j^* denotes D_j evaluated using \hat{F}_θ^*
 - 7: Update (a, b, α) via gradient descent (ADAM) to minimize the KS-cal
 - 8: **end while**
 - 9: Apply final calibrated transformation to the test set using optimized (a, b, α)
 - 10: **Output:** Calibrated CDF \hat{F}_θ^*
-

- No

- surrogate loss
 - additional nonparametric estimator
 - quantile estimation
 - sampling procedure

- Yes

- intuitive and easy to implement
 - preserve time-dependent C-index

Result

Table 1: Summary of pairwise comparisons between post-processing methods. The table shows the number of cases where KSP outperforms its counterpart, is outperformed, or yields a tie. Numbers in parentheses indicate statistically significant differences based on a one-sided t -test at the 0.05 level.

Method	C-index	S-cal(20)	D-cal(20)	KS-cal	KM-cal	IBS
KSP	20 (12)	46 (45)	46 (43)	47 (45)	47 (35)	44 (25)
Non-calibrated	18 (1)	13 (7)	14 (6)	13 (5)	13 (10)	16 (1)
Ties	22	1	0	0	0	0
KSP	13 (2)	36 (29)	48 (45)	51 (42)	37 (32)	42 (25)
CSD	34 (2)	24 (19)	12 (10)	9 (8)	23 (19)	18 (10)
Ties	13	0	0	0	0	0
KSP	21 (0)	32 (21)	46 (39)	44 (29)	45 (36)	38 (9)
CSD-iPOT	25 (1)	28 (19)	14 (13)	16 (11)	15 (10)	22 (10)
Ties	14	0	0	0	0	0

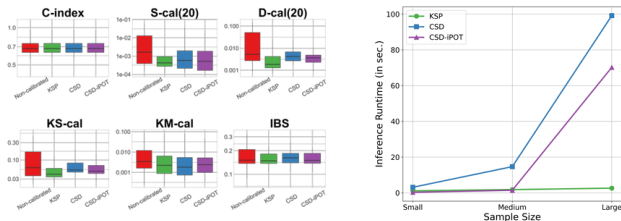


Figure 2: Boxplots of metric values (left) and inference runtime by sample size (right), aggregated across all datasets and models.

Conclusion

- Strength
 - Capture local discrepancies more than quantile-based methods
 - Scalable
- Weakness
 - Sensitive to discretized models
 - Less robust to tied times
- Future research
 - Extend KSP to conditional calibration