

Deep Copula-Based Survival Analysis for Dependent Censoring with Identifiability Guarantees

WeiJia Zhang^{1*} Chun Kai Ling² Xuanhui Zhang³

¹ School of Information and Physical Sciences, The University of Newcastle, Australia

² Carnegie Mellon University, USA

³ School of Information Management, Nanjing University, China

weijia.zhang.xh@gmail.com, chunkail@cs.cmu.edu, xuanhui@nju.edu.cn

Abstract

Censoring is the central problem in survival analysis where either the time-to-event (for instance, death), or the time-to-censoring (such as loss of follow-up) is observed for each sample. The majority of existing machine learning-based survival analysis methods assume that survival is conditionally independent of censoring given a set of covariates; an assumption that cannot be verified since only marginal distributions are available from the data. The existence of dependent censoring, along with the inherent bias in current estimators has been demonstrated in a variety of applications, accentuating the need for a more nuanced approach. However, existing methods that adjust for dependent censoring require practitioners to specify the ground truth copula. This requirement poses a significant challenge for practical applications, as model misspecification can lead to substantial bias. In this work, we propose a flexible deep learning-based survival analysis method that simultaneously accommodates for dependent censoring and eliminates the requirement for specifying the ground truth copula. We discuss the identifiability of our model under a broad family of copulas and survival distributions. Experiment results from a wide range of datasets demonstrate that our approach successfully discerns the underlying dependency structure and significantly reduces survival estimation bias when compared to existing methods.

1 Introduction

Survival analysis is a branch of statistical methods that focuses on modeling the time it takes for certain events to occur, with seminal work tracing back to mid-twentieth century such as the Kaplan-Meier estimator (Kaplan and Meier 1958) and Cox partial likelihood (Cox 1972). Survival analysis has been widely applied in many disciplines, with applications in healthcare such as epidemiology (Selvin 2008), clinical trials (Emmerson and Brown 2021), and personalized medicine (Zhang et al. 2017), as well as equipment failure time analysis (Voronov, Frisk, and Krysanter 2018).

The most prominent challenge of survival analysis is the existence of *censoring*, which occurs when the event time of a sample is not fully observed. Censoring is ubiquitous in clinical trials because a participant has the right to withdraw (Ondrusek et al. 1998). For example, in clinical follow-up

study designed to evaluate the effect of radiology, the period of relapse-free survival can only be observed if a participant exhibits the expected symptoms during the span of the follow-up. However, the true time-to-event time remains unobservable if the participant withdraws prematurely or the study concludes prior to the cancer relapse. This inherent uncertainty requires algorithms that account for censoring. Neglecting censored observations can result in loss of efficiency and estimation bias unless the observations are missing completely at random (Leung, Elashoff, and Afifi 1997).

A common assumption underpinning most machine learning and statistical survival analysis stipulates that censoring and survival are *conditionally independent* given the observed covariates (Wang, Li, and Reddy 2019). This assumption allows censored observations to be utilized by simultaneously maximizing the log-likelihood for both censored and uncensored samples. Unfortunately, as only marginal distributions of event and censoring are available from data, the independent-censoring assumption cannot be verified from observational data (Tsiatis 1975) in the similar sense that the unconfoundedness assumption is unverifiable in causal inference (Rubin 1974).

In many real world applications, censoring mechanisms are in fact *dependent* (Kaplan and Meier 1958; Leung, Elashoff, and Afifi 1997; Templeton, Amir, and Tannock 2020). For example, participants in clinical trials often prematurely remove themselves from the trial if they find the drug to be ineffective or experience adverse effects (Scharfstein, Rotnitzky, and Robins 1999). Similarly, in observational epidemiology, patients with a more advanced disease might be more likely to miss their follow-up clinic visits (Howe et al. 2010). Ignoring this dependence can result in biased survival estimations (Kleinbaum 2012).

One recent approach to account for the dependencies is to assume the parametric form of the dependencies and learn the correlation parameter (Emura et al. 2018; Deresa and Van Keilegom 2021; Gharari et al. 2023). These methods utilize *copulas*, which are powerful statistical tools which model dependencies between random variables in isolation from their marginals. That is, if the copula between observed and censored times is known, then survival marginals may be unbiasedly estimated even in the presence of censored data. These methods face two significant and interrelated challenges. Firstly, specifying the parametric family of

*Corresponding author.

the underlying copula is inherently difficult, as practitioners may lack experience or prior information to make an informed choice. Secondly, misspecified copulas will exacerbate model bias, leading to incorrect inference and misleading results. More recent methods by Deresa and Van Keilegom alleviate these issues by *learning the association parameters* of the copula. Nonetheless, the problem of misspecifying the copula’s parametric form still exists.

In this paper, we tackle dependent censoring by eliminating the requirement for a pre-specified copula. We demonstrate that the copula characterizing the dependency structure can be identified under reasonably mild conditions and propose learning them using deep neural networks which are trained end-to-end alongside the marginals. Specifically:

- We propose the Deep Copula Survival (DCSURVIVAL) framework¹, a deep copula-based survival analysis method that addresses *dependent censoring without requiring users to specify the parametric form of the ground truth copula*.
- We discuss the theoretical properties of our framework, demonstrating that under mild assumptions *identifiability* is attainable with common parametric survival marginals and the Archimedean copula family.
- We evaluate our method on a variety of datasets, demonstrating that DCSURVIVAL successfully learns the underlying copula and *significantly reduces bias* in survival predictions when compared to existing state-of-the-art.

2 Survival, Censoring and Copula

We are working with a survival dataset \mathcal{D} with the i -th sample denoted by $\mathcal{D}_i = (\mathbf{x}_i, t_i, \delta_i)$, where $\mathbf{x}_i \in \mathbb{R}^n$ is the n -dimensional covariates; $t_i \in \mathbb{R}^+$ is the observed time; and $\delta_i \in \{0, 1\}$ is the event indicator. We focus on the common scenario of *right censoring*, where $t_i = \min(T_i, U_i)$ with $T_i, U_i \in \mathbb{R}^+$ denoting the latent event and censoring times respectively. We have $\delta_i = 1$ when the event time is observed, while $\delta_i = 0$ when only the censored time is observed. We omit the subscript i when the context is clear. Under this model, the likelihood of a survival data point (\mathbf{x}, t, δ) under right-censoring is (Emura and Chen 2018)

$$\mathcal{L} = \Pr(T = t, U > t | \mathbf{x})^\delta \Pr(T > t, U = t | \mathbf{x})^{1-\delta} \quad (1)$$

We denote the marginal distributions for event and censoring time by $S_{T|X}(t|\mathbf{x}) = \Pr(T > t | \mathbf{x})$ and $S_{U|X}(u|\mathbf{x}) = \Pr(U > u | \mathbf{x})$, with density functions $f_{T|X}(t|\mathbf{x}) = -\partial S_{T|X}(t|\mathbf{x}) / \partial t$ and $f_{U|X}(t|\mathbf{x}) = -\partial S_{U|X}(t|\mathbf{x}) / \partial t$.

Independent Censoring. Most existing models assume that survival and censoring are independent, i.e., $T_i \perp\!\!\!\perp U_i$; or conditionally independent given the covariates, i.e., $T_i \perp\!\!\!\perp U_i | \mathbf{x}_i$. The likelihood function in Equation 1 simplifies to

$$\begin{aligned} \mathcal{L}_{\text{indep}} &= f_{T|X}(t|\mathbf{x}) S_{U|X}(t|\mathbf{x})^\delta \cdot f_{U|X}(t|\mathbf{x}) S_{T|X}(t|\mathbf{x})^{1-\delta} \\ &\propto f_{T|X}(t|\mathbf{x})^\delta \cdot S_{T|X}(t|\mathbf{x})^{1-\delta}. \end{aligned} \quad (2)$$

The density $f_{U|X}(t|\mathbf{x})$ and survival function $S_{U|X}(t|\mathbf{x})$ for the censoring distribution are often omitted during optimization as they are non-informative to the event densities and survival distributions (Kleinbaum 2012). However, in observational studies, the censoring mechanism is not only often unknown to researchers, but also unidentifiable solely based on observational data (Tsiatis 1975). This motivates methods to address dependent censoring without relying on specific or pre-determined assumptions about the censoring mechanism. One method is to utilize *copulas*.

Copulas. Copulas model dependencies between random variables in isolation from their marginal distributions. Loosely speaking, $\mathcal{C}(u_1, \dots, u_d) : [0, 1]^d \rightarrow [0, 1]$ is a d -dimensional copula if it is a distribution function of a random variable with support $[0, 1]^d$ and uniform marginals.

Copula has found extensive applications thanks to Sklar’s theorem (Sklar 1959), which states that any d -dimensional continuous joint distribution can be *uniquely* expressed with d uniform marginals and a copula \mathcal{C} . More formally,

Theorem 1. (Sklar 1959). *Let F be a distribution function with margins F_1, \dots, F_d . There exists a d -dimensional copula \mathcal{C} such that for any $(x_1, \dots, x_d) \in \mathcal{R}^d$ we have $F(x_1, \dots, x_d) = \mathcal{C}(F(x_1), \dots, F(x_d))$. Furthermore, if the marginals F_1, \dots, F_d are continuous, \mathcal{C} is unique.*

In practice, *Archimedean* copulas such as Clayton, Frank, Gumbel, and Joe copulas are common. Archimedean copulas are defined based on 1-dimensional generator φ , where

$$\mathcal{C}(u_1, \dots, u_d) = \varphi(\varphi^{-1}(u_1) + \dots + \varphi^{-1}(u_d)). \quad (3)$$

Here $\varphi : [0, \infty] \rightarrow [0, 1]$ is d -monotone, i.e., $(-1)^k \varphi^{(k)}(u) \geq 0$ for $k \leq d$ and $u \geq 0$. We say that φ is *completely monotone* if $(-1)^k \varphi^{(k)}(u) \geq 0$ for all $k \geq 0$.

Dependent Censoring via Copula. Dependent censoring arises when unobserved confounders affect both survival and censoring times, leading to dependencies that must be accounted for when evaluating joint likelihoods in Equation 1. This is similar to confounding in causal inference (Figure 1): in dependent censoring, we never simultaneously observe the censoring and survival times for a subject; while in causal inference, we never observe the factual and counterfactual outcomes at the same time (Pearl 2009).

When survival and censoring times are dependent, applying Sklar’s theorem yields the more general expression

$$\Pr(T > t, U > u | \mathbf{x}) = \mathcal{C}(S_{T|X}(t|\mathbf{x}), S_{U|X}(u|\mathbf{x})), \quad (4)$$

¹<https://github.com/WeijiaZhang24/DCSurvival>

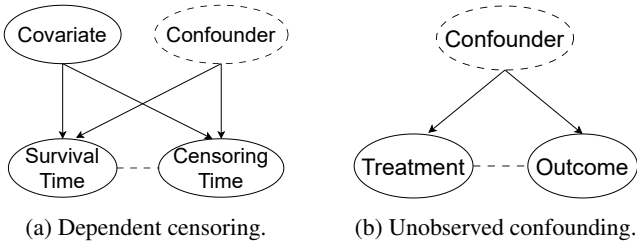


Figure 1: Illustrations of (a) dependent censoring in survival analysis and (b) unobserved confounding in causal inference. Solid/dashed nodes denote observed/hidden variables, respectively. The dashed lines between survival/censoring time and treatment/outcome indicate that estimation and evaluation in survival analysis with dependent censoring face similar challenges as in causal inference.

which when combined with Equation 1 yields the likelihood

$$\begin{aligned}
\mathcal{L}_{\text{dep}} = & \left\{ -\frac{\partial}{\partial u_1} \Pr(T > u_1, U > t | X = \mathbf{x}) \Big|_{u_1=t} \right\}^{\delta} \\
& \cdot \left\{ -\frac{\partial}{\partial u_2} \Pr(T > t, U > u_2 | X = \mathbf{x}) \Big|_{u_2=t} \right\}^{1-\delta} \\
= & \left\{ f_{T|X}(t|\mathbf{x}) \frac{\partial}{\partial u_1} \mathcal{C}(u_1, u_2) \Big|_{\substack{u_1=S_{T|X}(t|\mathbf{x}) \\ u_2=S_{U|X}(t|\mathbf{x})}} \right\}^{\delta} \\
& \cdot \left\{ f_{U|X}(t|\mathbf{x}) \frac{\partial}{\partial u_2} \mathcal{C}(u_1, u_2) \Big|_{\substack{u_1=S_{T|X}(t|\mathbf{x}) \\ u_2=S_{U|X}(t|\mathbf{x})}} \right\}^{1-\delta}.
\end{aligned} \tag{5}$$

In this paper, we make the mild assumption that \mathcal{C} does not depend on \mathbf{x} . We see that \mathcal{L}_{dep} and $\mathcal{L}_{\text{indep}}$ are equivalent only when \mathcal{C} is the independence copula $\mathcal{C}(u_1, u_2) = u_1 u_2$, which corresponds to the case where T and U are conditionally independent. For all other copulas, $\mathcal{L}_{\text{indep}}$ is biased because the dependency between T and U and the censoring marginals are not ignorable.

To the best of our knowledge, all existing copula-based survival models require practitioners to *specify the family of ground truth copula*. As most widely-used bivariate copula have closed-form partial derivatives, \mathcal{L}_{dep} can then be optimized given the parameterization of the survival distributions and the assumed copula. However, correctly defining the true copula presents an inherent challenge, placing practitioners in a difficult position as any misspecification can substantially amplify bias.

3 End-to-end Survival Analysis via Copula

We now introduce DCSURVIVAL, a framework seeking to learn both copula and survival distributions directly from right-censored data. In particular, we jointly maximize the likelihood in Equation 5, fitting parameters of marginal distributions of $T|X, U|X$ and the parameters of the copula \mathcal{C} between them. DCSURVIVAL optimizes over copula belonging to the Archimedean family; crucially, this includes

the bivariate independence copula $\mathcal{C}(u_1, u_2) = u_1 u_2$ and hence subsumes the independent-censoring assumption.

DCSURVIVAL comprises two components, each based on the terms in Equation 5. The first is the Archimedean copula \mathcal{C} between T and U , the second comprises the density functions for event and censoring times $f_{T|X}$ and $f_{U|X}$. At a high level, suppose $\mathcal{C}(u_1, u_2)$, $f_{T|X}$, and $f_{U|X}$ are parameterized by $\alpha = (\theta, \theta_U, \theta_T)$ respectively, and that \mathcal{C} is specified such that the partial derivatives $\partial \mathcal{C}(u_1, u_2) / \partial u_1$ and $\partial \mathcal{C}(u_1, u_2) / \partial u_2$ alongside their derivatives with respect to θ may be computed for all u_1, u_2 . Then, one can compute \mathcal{L}_{dep} , and by applying the chain and product rules obtain the required gradients $\partial \mathcal{L}_{\text{dep}} / \partial \alpha$ needed to optimize α in an end-to-end fashion via gradient descent.

Using this framework, one can restrict $U|X, T|X$ and \mathcal{C} to be “textbook” distributions with few parameters to be estimated. This however leads to problems of model misspecification. DCSURVIVAL allows for the use of neural networks to model both copulas and margins, as long as they are “fully differentiable” in the manner described above. We describe both components and how they relate to our problem of dependent censoring. For brevity, we defer specific architectural details to supplemental material.

Archimedean Copula for Event and Censoring Times.

Kimberling (1974) showed that the generator φ being completely monotone is a necessary and sufficient condition for all *two-dimensional* Archimedean copulas. Furthermore, utilizing the Bernstein-Widder characterization theorem (Widder 2010), we know that every completely monotone function can be characterized using a mixture of negative exponential functions. Formally, we have:

Theorem 2. (Bernstein-Widder) *A generator φ is completely monotone if and only if φ is the Laplace transform of a positive random variable M , i.e., $\varphi(u) = \mathbb{E}_M[\exp(-uM)]$ and $P(M > 0) = 1$.*

Combining Theorem 2 with the results of Kimberling (1974) implies that any completely monotone generator φ can be approximated using a *finite* sum of negative exponentials (Koyama 2023). This includes all bivariate Archimedean copulas thanks to the necessary condition.

As such, we learn the generator φ (which implicitly defines the copula via Equation 3 for $d = 2$) via the method by Ling, Fang, and Kolter (2020), which uses neural networks to parameterize φ with a large but finite mixture of negative exponentials. Roughly speaking, the neural network $\varphi_{nn} : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ resembles a multilayer fully-connected network with a 1-dimensional input, where each neuron’s output is a convex combination of negative exponentials (of the input the network) with different rates. These are mixed with other neurons in the same layer (with some bias included) based on the network weights to give the outputs of the neurons of the next layer. As Ling, Fang, and Kolter (2020) show, this rich architecture allows φ_{nn} to satisfy Theorem 2 and be a convex combination of an exponential (in the size of the network) number of negative exponentials.

Parameterizing an Archimedean copula via its generator has many benefits, chief of which is that we are able to compute $\frac{\partial}{\partial u_1} \mathcal{C}(u_1, u_2)$ and $\frac{\partial}{\partial u_2} \mathcal{C}(u_1, u_2)$ evaluated at $u_1 =$

$S_{T|X}(t|\mathbf{x})$ and $u_2 = S_{U|X}(t|\mathbf{x})$ such that their derivatives with respect to the network parameters θ may in turn be used for gradient descent. In practice, this process is simplified via the use of automatic differentiation tools like Pytorch (Paszke et al. 2017). We note that computing these quantities will involve evaluating Equation 3, which necessitates estimating the inverse of φ_{nn} . This can be done efficiently via Newton’s method since φ_{nn} is one dimensional, and we refer to reader to Ling, Fang, and Kolter (2020) for details.

Remark. The network architecture of φ_{nn} is not restricted to bivariate Archimedean copula and can be applied to survival analysis with competing risks. However, a key difference is that the completely monotone condition in Theorem 2 is not necessary for the generators for Archimedean copulas. However, all *extendible* Archimedean copulas can still be expressed with completely monotone generators when $d > 2$ (McNeil and Nešlehová 2009). In this work, we will focus on dependent censoring in single-risk and leave research on competing risks for future endeavours.

Survival and Censoring Marginals. When it is reasonable to assume the parametric form of the survival and censoring distributions, the corresponding survival function $S_{T|X}(t|\mathbf{x})$ and the event density function $f_{T|X}(t|\mathbf{x})$ can be directly plugged into Equation 5 for evaluation. For example, CoxPH model with Weibull marginals (Gharari et al. 2023) or log-normal marginals can be used. In this case, identifiability is guaranteed: we discuss this in Section 4.

When the parametric form for survival marginals is unknown, one can represent $S_{T|X}$ and $S_{U|X}$ using the monotonic neural density estimators (NDE) proposed by Chilinski and Silva (2020). The NDE network consists of a fully-connected covariate network for learning the representation of X , alongside a monotonic network parameterized with non-negative weights to ensure decreasingness with respect to t . Specifically, the output of the covariate network is concatenated with t and fed to the layers with non-negative weights. The final layer of the monotonic network outputs a scalar with sigmoid activation yielding the survival function. Correspondingly, the density function may be expressed by $f_{T|X}(t|\mathbf{x}) = -\partial S_{T|X}(t|\mathbf{x})/\partial t$ and $f_{U|X}(t|\mathbf{x}) = -\partial S_{U|X}(t|\mathbf{x})/\partial t$, each computed via auto differentiation.

Once both marginals and copula are instantiated, their parameters are estimated by maximizing Equation 5 via stochastic gradient descent. Note that aside from survival distributions, DCSURVIVAL also estimates \mathcal{C} as an added bonus. This may be of independent interest to practitioners.

4 Model Identifiability

Unidentifiable models hold no guarantee that the true parameters are recovered even with an infinite number of datapoints. Tsiatis (1975) showed that the joint distribution of the survival time T and censoring time C is not identifiable under a completely non-parametric setting. Therefore, we discuss identifiability under the mild assumption that the dependency can be characterized by the Archimedean family of copulas. Furthermore, our approach also has the benefit of not making strict assumptions regarding the survival marginals. To our knowledge, existing previous work

(Gharari et al. 2023) is not identifiable and requires practitioners to specify the exact form of the ground truth copula.

Specifically, we consider identifiability for the following model with parameters $\alpha = (\theta, \theta_T, \theta_U)$. In this section, we make the dependence on parameters more explicit (while omitting covariates \mathbf{x} for brevity), yielding the following.

$$\begin{aligned} \Pr(T > t, U > u) &= \mathcal{C}_\theta(S_{T,\theta_T}(t), S_{U,\theta_U}(u)), \\ f_{y,\delta=1,\alpha}(y) &= f_{T,\theta_T}(y) \left. \frac{\partial}{\partial u_1} \mathcal{C}_\theta(u_1, u_2) \right|_{\substack{u_1=S_{T,\theta_T}(t) \\ u_2=S_{U,\theta_U}(t)}}, \\ f_{y,\delta=0,\alpha}(y) &= f_{U,\theta_U}(y) \left. \frac{\partial}{\partial u_2} \mathcal{C}_\theta(u_1, u_2) \right|_{\substack{u_1=S_{T,\theta_T}(t) \\ u_2=S_{U,\theta_U}(t)}}, \end{aligned} \quad (6)$$

where $f_{y,\delta=1,\alpha}(y)$ and $f_{y,\delta=0,\alpha}(y)$ are the densities of event and censored observations.

We start with the sufficient condition for identifying the model in Equation 6. Specifically, we will show that the parameters $\alpha = (\theta, \theta_U, \theta_T)$ can be uniquely determined from the observed dataset $\mathcal{D} = \{(\mathbf{x}_i, t_i, \delta_i)\}$. In other words, if $f_{y,\delta,\alpha_1}(y) = f_{y,\delta,\alpha_2}(y)$ for all y , then $\alpha_1 = \alpha_2$.

Theorem 3. (Czado and Van Keilegom 2023) Suppose the following conditions are satisfied: (C1) For all $\theta_{T_1}, \theta_{T_2} \in \Theta_T$ and $\theta_{U_1}, \theta_{U_2} \in \Theta_U$, we have

$$\begin{aligned} \lim_{t \rightarrow 0} \frac{f_{T,\theta_{T_1}}(t)}{f_{T,\theta_{T_2}}(t)} &= 1 \iff \theta_{T_1} = \theta_{T_2}, \text{ and} \\ \lim_{t \rightarrow 0} \frac{f_{U,\theta_{U_1}}(t)}{f_{U,\theta_{U_2}}(t)} &= 1 \iff \theta_{U_1} = \theta_{U_2}. \end{aligned} \quad (7)$$

(C2) For all $(\theta, \theta_T, \theta_U) \in \Theta \times \Theta_T \times \Theta_U$, we have

$$\begin{aligned} \lim_{t \rightarrow 0} \left. \frac{\partial}{\partial u_1} \mathcal{C}_\theta(u_1, u_2) \right|_{\substack{u_1=S_{T,\theta_T}(t) \\ u_2=S_{U,\theta_U}(t)}} &= 1 \\ \lim_{t \rightarrow 0} \left. \frac{\partial}{\partial u_2} \mathcal{C}_\theta(u_1, u_2) \right|_{\substack{u_1=S_{T,\theta_T}(t) \\ u_2=S_{U,\theta_U}(t)}} &= 1 \end{aligned} \quad (8)$$

Then, the model defined in Equation 6 is identified.

Proof Sketch. The proof consists of two main steps. First, identification for the marginal distributions of the event and censoring times can be shown by tying the densities of observed event/censoring times to the true density of event and censoring distributions utilizing Condition (C2). In a sense, this bears similarity to showing that the observed outcomes equal to the potential outcomes by expectation in causal inference. Second, the copula parameters can be identified by leveraging the unique existence of copulas. \square

Theorem 4. (Czado and Van Keilegom 2023) Condition (C1) is satisfied by the families of Weibull, log-normal, log-logistic, and log-Student densities.

Theorem 3 and 4 are attributed to (Czado and Van Keilegom 2023). As it turns out, Archimedean copulas represented by φ_{nn} have the attractive property of satisfying the condition (C2).

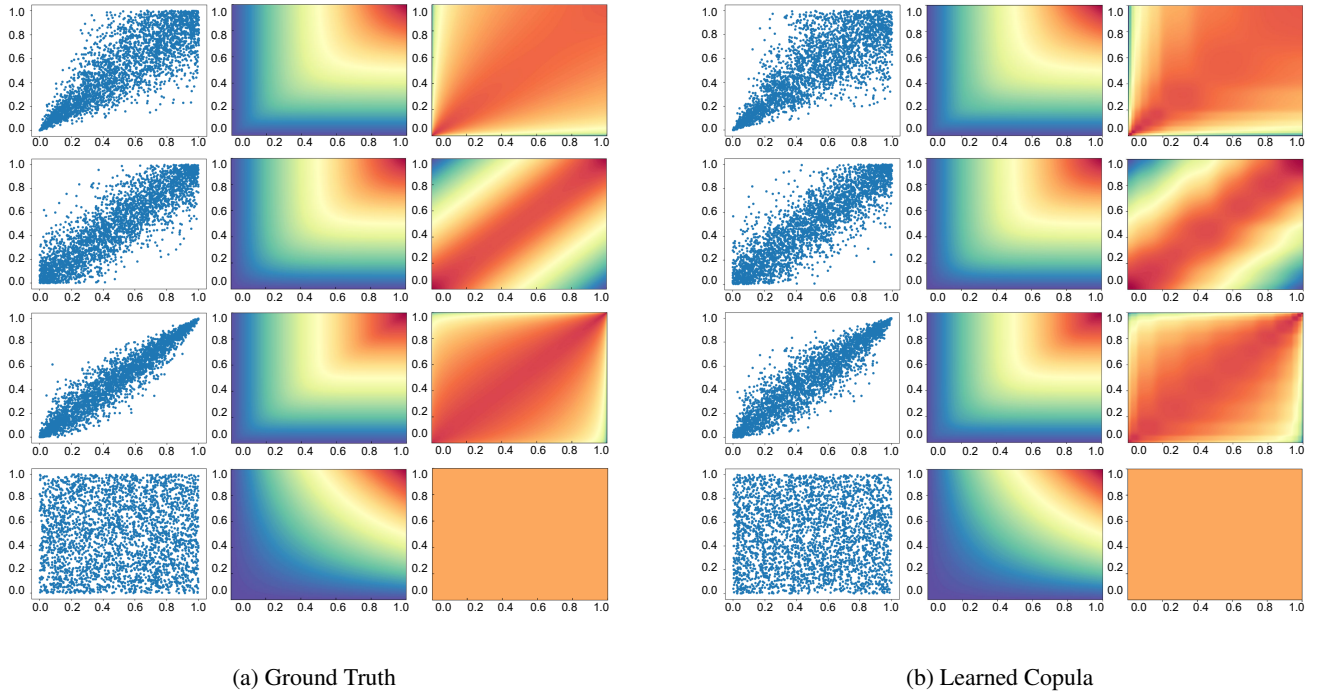


Figure 2: Top to bottom: learning Clayton, Frank, Gumbel, and Independence copulas using DCSURVIVAL from the Linear-Rsksk dataset. Left to right: ground truth versus learned copula for the (i) scatter sample plots; (ii) joint cumulative distribution plots, (iii) log probability density plots. These figures are best viewed in colour.

Lemma 5. *Condition (C2) of Theorem 3 is satisfied if φ is differentiable on $(0, \infty)$ and $\lim_{u \rightarrow 0} \varphi'(u) < 0$.*

Theorem 6. *Archimedean copulas represented by φ_{nn} satisfy condition (C2) in Theorem 3.*

Therefore, DCSURVIVAL is identifiable when (i) \mathcal{C} is defined using φ_{nn} and (ii) T and U have margins belonging to the families in Theorem 4. However, it is worth noting that Theorem 6 does not necessarily apply to all Archimedean copulas, but rather those that can be represented by φ_{nn} . In particular, Czado and Van Keilegom (2023) have shown that not all Archimedean copulas are identifiable. When using NDE to represent the survival marginals, the sufficient condition in (C1) of Theorem 3 is also satisfied since two neural networks with the same parameters produce the same outputs. However, the necessary condition is *not* guaranteed. Although the NDE instantiation does not strictly adhere to our identifiability condition, it compensates with superior flexibility, requiring neither proportional hazard nor survival marginal assumptions. Empirically, we find that DCSURVIVAL performs significantly better than state-of-the-art survival methods when dependent censoring is present, and achieves competitive performance under independent censoring, further highlighting its robustness and versatility.

5 Experiments

Evaluating dependent censoring methods is challenging due to the absence of a known censoring mechanism in observational datasets, akin to assessing treatment effect estimation in causal inference without ground truth causal effects

(Parikh et al. 2022). To address this, we employ (i) synthetic datasets with defined dependency structures, (ii) semi-synthetic datasets using real-world covariates for censoring time simulation, and (iii) real-world datasets using metrics that do not need ground truth, as commonly done in causal effect studies (Zhang, Li, and Liu 2021).

Experiments are conducted with one NVIDIA RTX4090 GPU. We utilize Pytorch (Paszke et al. 2017) for implementing all neural networks and automatic differentiation. Tensors are computed with double precision (fp64) as the inversion of φ mandates numerical precision. When using Newton’s method to compute the inverse φ_{nn}^{-1} , we terminate when the error is less than 1×10^{-12} . For all our experiments we set φ_{nn} with $L = 2$ and $H_1 = H_2 = 10$, i.e., the copula representation contains two hidden layers with each of width 10. The network is small but sufficiently expressive for the dependency structure since the generator φ_{nn} is only 1-dimensional. We use AdamW (Loshchilov and Hutter 2019) for optimization and use 50%/30%/20% training/validation/test splits. We used validation samples for early stopping based on the validation log-likelihood. No further hyperparameter tuning was performed. We provide our code and other details in the supplementary material.

Synthetic Datasets Following the approach of Gharari et al. (2023), we generate two synthetic datasets with the CoxPH model where the event and censoring risks are specified by Weibull marginal distributions (Bender, Augustin, and Blettner 2005). Specifically, we sample two random variables from known copula and then apply inverse trans-

form sampling to generate the event and censoring times T_i and U_i according to their Weibull distribution parameters. The observed t_i is the minimum of the two. We experiment with four Archimedean copulas, including the Clayton, Frank, Gumbel and Independence copulas and sample using copula-specific methods (Scherer and Mai 2012). We generate two datasets, *Linear-Risk* and *Nonlinear-Risk*, which correspond to cases where the Weibull hazards are linear and non-linear functions of covariates (which are in turn drawn from $\mathcal{U}_{[0,1]}^{10}$). Further details of the data generating procedure are provided in the supplementary materials.

Semi-Synthetic Datasets We use two semi-synthetic datasets based on the *STEEL* (V E, Shin, and Cho 2020) which contains 35,040 samples with 9 covariates, and the *Airfoil* (Thomas Brooks 1989) datasets which includes 1503 samples with 6 covariates. We induce censoring following a setting similar to those described in (Gharari et al. 2023). Briefly speaking, we use the dependent variable as the event time, and conditionally sample the censoring time with a copula. Contrary to the two synthetic datasets, the proportional hazard assumption is not maintained in the semi-synthetic datasets. The semi-synthetic dataset generation method is similar to widely-used causal effect estimation benchmarks (Hill 2011; Dorie et al. 2019).

Real-World Datasets We use two real-world datasets. The *SEER* dataset is from the Surveillance, Epidemiology and End Results database (Howlander et al. 2010). Following Czado and Van Keilegom (2023), we use the monthly survival time of patients with localized pancreas cancer diagnosed between 2000 and 2015, and the event of interest is death caused by pancreas cancer. Patients that are alive or have died because of other cancers are considered as censored. There are 15 covariates that measure demographic and pathology features of 11,600 participants. However, as many common risk factors of diseases are not measured, the censoring and event time is likely to be correlated when patients are censored due to death caused by other diseases.

GBSG2 is from the German Breast Cancer Study Group (Schumacher et al. 1994), which contains samples obtained from an observational study of 686 women with 8 covariates measuring the pathology characteristics of the participants. The event of interest is the relapse-free survival time, while the participants are censored when they pass away. Both are likely positively correlated due to the correlations between cancer recurrence and patient death. Further details of all datasets are provided in the supplementary material.

Identifying Censoring Dependency Structure

We first empirically validate our identifiability results by instantiating *DCSURVIVAL* with parametric survival marginals that satisfy (C1). Figure 2 shows the results of *DCSURVIVAL* for learning the copula from the *Linear-Risk* dataset by contrasting the sample scatter plots, joint cumulative distribution functions, and joint probability density functions of the ground truth copula with those learned by *DCSURVIVAL*. From Figure 2, we can see that *DCSURVIVAL* is able to learn the dependency structure

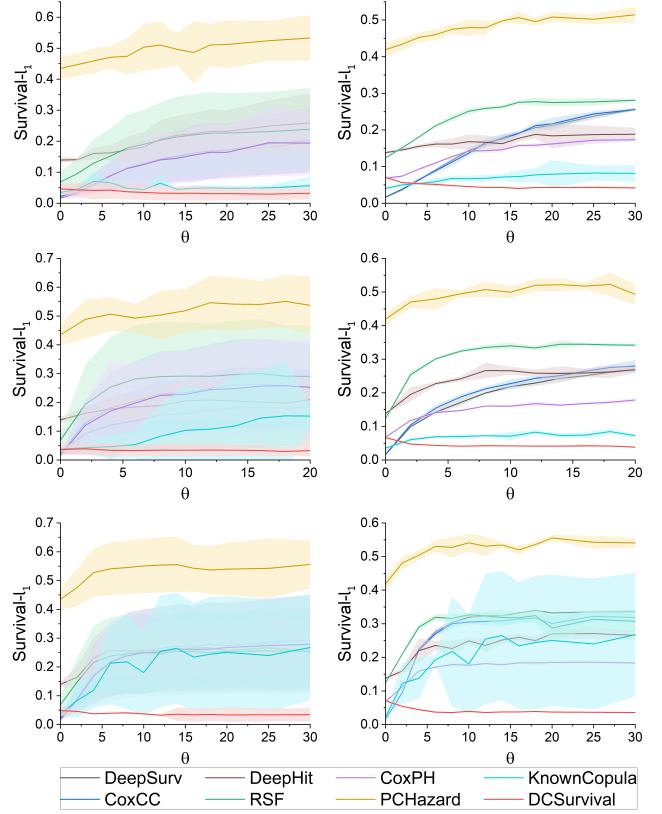


Figure 3: Top to bottom (row): survival prediction biases of compared algorithms on varying censoring dependency governed by Frank, Clayton, and Gumbel copulas. Left to right: *LINEAR-RISK* and *NONLINEAR-RISK* results. The lines represent the $\text{Survival-}l_1$ means and the shaded areas are the standard deviations. Best viewed in colour.

specified by Clayton, Frank, and Gumbel copulas from censored samples, and the contours of the log-likelihood almost exactly match the ground truth. Additionally, *DCSURVIVAL* also correctly recovers the independent censoring mechanism when the ground truth is specified by the Independence copula. These empirical results corroborate the identifiability result in Theorem 3, demonstrating that dependency structure can be learned from right-censored data. Experiments on *Nonlinear-Risk* dataset show similar good results which are provided in the supplementary material.

Reducing Survival Estimation Bias

We now assess the survival bias mitigation capabilities of *DCSURVIVAL* with the end-to-end NDE instantiation. Under dependent censoring, frequently used survival metrics such as C-indices (Harrell 1982) and Brier scores (Brier 1950) are not proper scoring rules, i.e., the highest score is not achieved by the true distribution (Gharari et al. 2023). Therefore, when the ground truths are known, we compare them with our estimated survival distributions using the

	STEEL	Airfoil
COXPH	0.246 \pm 0.008	0.235 \pm 0.054
DEEPSURV	0.234 \pm 0.007	0.265 \pm 0.031
COXCC	0.265 \pm 0.025	0.264 \pm 0.100
RSF	0.249 \pm 0.009	0.474 \pm 0.034
DEEPHIT	0.157 \pm 0.020	0.259 \pm 0.055
PCHAZARD	0.212 \pm 0.039	0.200 \pm 0.038
KNOWNCOPULA	0.105 \pm 0.013	0.215 \pm 0.060
DCSURVIVAL	0.092 \pm 0.015	0.176 \pm 0.047

Table 1: The Survival- l_1 (mean \pm std) metrics of compared methods on test samples. The results reported are the averages obtained from repeating each experiment for 10 times.

Survival- l_1 metric (Gharari et al. 2023):

$$\mathcal{C}_{\text{Survival-}l_1} = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{|\mathcal{D}|_{t_i^{\max}}} \int_0^{\infty} |S(t|\mathbf{x}_i) - \hat{S}(t|\mathbf{x}_i)| dt, \quad (9)$$

where S and \hat{S} are the ground truth and estimated survival distributions, respectively. For real-world datasets, we utilize *calibration plots* (Niculescu-Mizil and Caruana 2005) for evaluation, which can be obtained from the event indicators without accessing the true survival marginals.

We evaluate DCSURVIVAL against a diverse set of survival analysis techniques, including traditional methods COXPH (Cox 1972) and Random Survival Forest (RSF) (Ishwaran et al. 2008), as well as deep learning-based approaches such as DEEPSURV (Katzman et al. 2018), DEEPHIT (Lee et al. 2018), and PCHAZARD (Kvamme and Borgan 2021). We also compare with a copula-based method that assumes dependent censoring (Gharari et al. 2023), henceforth referred to as KNOWNCOPULA, as it requires known ground truth copula. Notably, KNOWNCOPULA also requires users to specify event and censoring marginals.

Results from Figure 3 and Table 1 show that DCSURVIVAL consistently achieves the best performances under dependent censoring. Observe that the performance of KNOWNCOPULA deteriorates as the dependency increases. This is possibly caused by the fact that explicit optimizing for θ involves calculating the exponentials of the inverse of θ , posing substantial numerical challenges. For independent censoring as shown in Figure 3, DCSURVIVAL also performs competitively by consistently outperforming RSF, DEEPHIT, and PCHAZARD. Importantly, algorithms surpassing DCSURVIVAL under independent censoring are restricted to the proportional hazard assumption which aligns with the data generation procedure.

In Figure 4, calibration plots for the SEER and GBSG2 datasets demonstrate that DCSURVIVAL is the closest to the ideal calibration represented by the black dashed line. Calibration plot assesses whether an algorithm underestimates/overestimates event risks by comparing the average predicted event risk with the overall event rate (Calster et al. 2019). Other compared methods systematically underestimate the event probability, as evidenced by their plots residing substantially above the ideal calibration. For KNOWNCOPULA, because the ground truth copula is unknown, we

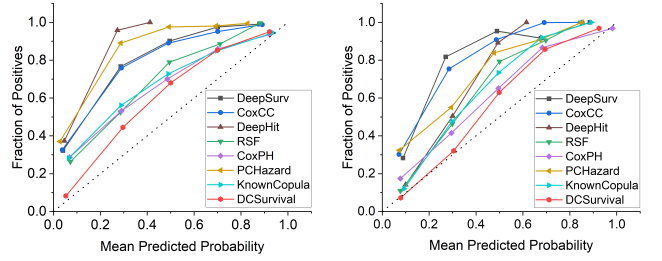


Figure 4: Calibration plots of on test samples of the SEER (left) and GBSG2 (right) datasets. The plots of better-calibrated algorithms are closer to perfect calibration (dashed black line). Best viewed in colour.

use a convex combination of Frank and Clayton copulas, and simultaneously optimizing for the mixing weight and copula parameters as discussed in (Gharari et al. 2023). However, this approach does not yield good performance, possibly due to misspecification of both the copula and survival marginals. These observations emphasize the value of data-driven methods like DCSURVIVAL, which jointly learns copula and survival marginals.

6 Related Work

Independent Censoring Methods. Most statistical survival models assume independent censoring. Notably, the partial likelihood function introduced in the Cox Proportional Hazard (CoxPH) model (Cox 1972) explicitly relies on the independent censoring assumption (Jackson et al. 2014). Furthermore, classical models such as the Kaplan-Meier estimator (Kaplan and Meier 1958) and Nelson-Aalen estimator (Nelsen 1998) also assume independent censoring. For machine learning-based methods, Random Survival Forest (Ishwaran et al. 2008) ensembles survival trees build with the log-rank test splitting criterion which assumes independent censoring. Neural network-based approaches, such as Fraggi-Simon Net (Fraggi and Simon 1995), DeepSurv (Katzman et al. 2018), and CoxCC (Kvamme, Ørnulf Borgan, and Scheel 2019) adhere to the proportional hazard assumption of the CoxPH model, and thus also mandate independence. DeepHit (Lee et al. 2018) discretizes the event time horizon and accommodates competing risks with event-specific neural networks. Recently, several tailored neural models such as differential equation networks (Tang et al. 2022) and monotonic networks (Rindt et al. 2022) have also been proposed for survival data. Unfortunately, they all assume independent or conditional independent censoring.

Dependent Censoring Methods. Discussion on dependent censoring in survival analysis can be traced back to the seminal work of Tsiatis (1975) and Lagakos (1979). To account for dependent censoring, several copula-based approaches have been explored by assuming the parametric family of the ground truth copula is known. Specifically, Czado and Van Keilegom (2023); Deresa and Van Keilegom (2021, 2024) proposed statistical methods for parametric and semiparametric survival marginals with identifiability

ity guarantee for all parameters. Midtjord, Bin, and Huseby (2022) and Gharari et al. (2023) proposed to use boosting and neural network for maximizing the log-likelihood with pre-specified copulas, respectively. To the best of our knowledge, all existing copula-based approaches require practitioners to provide the ground truth copula, or at least its parametric form. Our proposed DCSURVIVAL algorithm relaxes this assumption by approximating the copula with neural networks and learning their parameters.

7 Conclusion

In this paper, we propose DCSURVIVAL, a deep copula-based survival algorithm that does not require user-specified ground truth copula. DCSURVIVAL is capable of survival modelling under a wide range of censoring dependencies described by the Archimedean family of copulas. Theoretically, we show that the parameters for both the copula and the survival distributions are identifiable under mild assumptions. Empirically, we demonstrate that DCSURVIVAL successfully recovers the censoring dependencies and significantly reduces survival estimation bias. Future work include moving beyond Archimedean copulas and developing neural survival marginal estimators that are identifiable.

Acknowledgments

The authors would like to thank Prof. Ingrid Van Keilegom for her discussion and constructive comments on this work.

Xuanhui Zhang is supported by the National Science Foundation of China (72204110).

References

- Bender, R.; Augustin, T.; and Blettner, M. 2005. Generating survival times to simulate Cox proportional hazards models. *Statistics in medicine*, 24(11): 1713–1723.
- Brier, G. W. 1950. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1): 1–3.
- Calster, B. V.; McLernon, D. J.; van Smeden, M.; Wynants, L.; and Steyerberg, E. W. 2019. Calibration: the Achilles heel of predictive analytics. *BMC Medicine*, 17(1).
- Chilinski, P.; and Silva, R. 2020. Neural likelihoods via cumulative distribution functions. In *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, 420–429. PMLR.
- Cox, D. R. 1972. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2): 187–202.
- Czado, C.; and Van Keilegom, I. 2023. Dependent censoring based on parametric copulas. *Biometrika*, 110(3): 721–738.
- Davidson-Pilon, C. 2023. lifelines, survival analysis in Python.
- Deresa, N. W.; and Van Keilegom, I. 2021. On semiparametric modelling, estimation and inference for survival data subject to dependent censoring. *Biometrika*, 108(4): 965–979.
- Deresa, N. W.; and Van Keilegom, I. 2024. Copula based cox proportional hazards models for dependent censoring. *Journal of the American Statistical Association*, 119: 1044–1054.
- Dorie, V.; Hill, J.; Shalit, U.; Scott, M.; and Cervone, D. 2019. Automated versus do-it-yourself methods for causal inference: lessons learned from a data analysis competition. *Statistical Science*, 34(1): 43–68.
- Emmerson, J.; and Brown, J. 2021. Understanding survival analysis in clinical trials. *Clinical Oncology*, 33(1): 12–14.
- Emura, T.; and Chen, Y.-H. 2018. *Analysis of survival data with dependent censoring: copula-based approaches*. Springer.
- Emura, T.; Chen, Y.-H.; Matsui, S.; and Rondeau, V. 2018. *Analysis of survival data with dependent censoring: copula-based approaches*. Springer.
- Faraggi, D.; and Simon, R. 1995. A neural network model for survival data. *Statistics in Medicine*, 14(1): 73–82.
- Gharari, A. H. F.; Cooper, M.; Greiner, R.; and Krishnan, R. G. 2023. Copula-based deep survival models for dependent censoring. In *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 216 of *Proceedings of Machine Learning Research*, 669–680. PMLR.
- Harrell, F. E. 1982. Evaluating the yield of medical tests. *JAMA: The Journal of the American Medical Association*, 247(18): 2543.
- Hill, J. L. 2011. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1): 217–240.
- Howe, C. J.; Cole, S. R.; Napravnik, S.; and Eron, J. J. 2010. Enrollment, retention, and visit attendance in the university of north carolina center for AIDS research HIV clinical cohort, 2001–2007. *AIDS Research and Human Retroviruses*, 26(8): 875–881.
- Howlader, N.; Ries, L. A. G.; Mariotto, A. B.; Reichman, M. E.; Ruhl, J.; and Cronin, K. A. 2010. Improved Estimates of Cancer-Specific Survival Rates From Population-Based Data. *JNCI: Journal of the National Cancer Institute*, 102(20): 1584–1598.
- Ishwaran, H.; Kogalur, U. B.; Blackstone, E. H.; and Lauer, M. S. 2008. Random survival forests. *The Annals of Applied Statistics*, 2(3).
- Jackson, D.; White, I. R.; Seaman, S.; Evans, H.; Baisley, K.; and Carpenter, J. 2014. Relaxing the independent censoring assumption in the Cox proportional hazards model using multiple imputation. *Statistics in Medicine*, 33(27): 4681–4694.
- Kaplan, E. L.; and Meier, P. 1958. Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53(282): 457–481.
- Katzman, J. L.; Shaham, U.; Cloninger, A.; Bates, J.; Jiang, T.; and Kluger, Y. 2018. DeepSurv: personalized treatment recommender system using a cox proportional hazards deep neural network. *BMC Medical Research Methodology*, 18(1).

- Kimberling, C. H. 1974. A probabilistic interpretation of complete monotonicity. *Aequationes Mathematicae*, 10(2-3): 152–164.
- Kleinbaum, D. G. 2012. *Survival analysis a self-learning text*. Springer.
- Koyama, Y. M. 2023. Exponential sum approximations of finite completely monotonic functions. *arXiv:2301.08931*.
- Kvamme, H.; and Borgan, Ø. 2021. Continuous and discrete-time survival prediction with neural networks. *Life-time Data Analysis*, 27(4): 710–736.
- Kvamme, H.; Ørnulf Borgan; and Scheel, I. 2019. Time-to-event prediction with neural networks and cox regression. *Journal of Machine Learning Research*, 20(129): 1–30.
- Lagakos, S. W. 1979. General right censoring and its impact on the analysis of survival data. *Biometrics*, 35: 139–156.
- Lee, C.; Zame, W.; Yoon, J.; and van der Schaar, M. 2018. Deephit: a deep learning approach to survival analysis with competing risks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32: 2314–2321.
- Leung, K.-M.; Elashoff, R. M.; and Afifi, A. A. 1997. Censoring issues in survival analysis. *Annual Review of Public Health*, 18(1): 83–104.
- Ling, C. K.; Fang, F.; and Kolter, J. Z. 2020. Deep archimedean copulas. In *Advances in Neural Information Processing Systems*, volume 33, 1535–1545.
- Loshchilov, I.; and Hutter, F. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- McNeil, A. J.; and Nešlehová, J. 2009. Multivariate archimedean copulas, d-monotone functions and ℓ_1 -norm symmetric distributions. *The Annals of Statistics*, 37(5B).
- Midtfjord, A. D.; Bin, R. D.; and Huseby, A. B. 2022. A copula-based boosting model for time-to-event prediction with dependent censoring. *arXiv:2210.04869*.
- Nelsen, R. B. 1998. *An introduction to copulas*. Springer. ISBN 9780387986234.
- Niculescu-Mizil, A.; and Caruana, R. 2005. Predicting good probabilities with supervised learning. In *Proceedings of the 22nd international conference on Machine learning - ICML '05*, 625–632.
- Ondrusek, N.; Abramovitch, R.; Pencharz, P.; and Koren, G. 1998. Empirical examination of the ability of children to consent to clinical research. *Journal of Medical Ethics*, 24(3): 158–165.
- Parikh, H.; Varjao, C.; Xu, L.; and Tchetgen, E. T. 2022. Validating causal inference methods. In Chaudhuri, K.; Jegelka, S.; Song, L.; Szepesvari, C.; Niu, G.; and Sabato, S., eds., *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 17346–17358. PMLR.
- Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic Differentiation in PyTorch. In *NIPS 2017 Workshop on Autodiff*.
- Pearl, J. 2009. *Causality*. Cambridge, England: Cambridge University Press.
- Pölsterl, S. 2020. Scikit-survival: a library for time-to-event analysis built on top of scikit-learn. *Journal of Machine Learning Research*, 21(212): 1–6.
- Rindt, D.; Hu, R.; Steinsaltz, D.; and Sejdinovic, D. 2022. Survival regression with proper scoring rules and monotonic neural networks. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, 1190–1205. PMLR.
- Rubin, D. B. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5): 688–701.
- Scharfstein, D. O.; Rotnitzky, A.; and Robins, J. M. 1999. Adjusting for nonignorable drop-out using semiparametric nonresponse models. *Journal of the American Statistical Association*, 94(448): 1096–1120.
- Scherer, M.; and Mai, J. 2012. *Simulating copulas: stochastic models, sampling algorithms, and applications*. Imperial College Press.
- Schumacher, M.; Bastert, G.; Bojar, H.; Hübner, K.; Olschewski, M.; Sauerbrei, W.; Schmoor, C.; Beyerle, C.; Neumann, R. L.; and Rauschecker, H. F. 1994. Randomized 2 x 2 trial evaluating hormonal treatment and the duration of chemotherapy in node-positive breast cancer patients. German breast cancer study group. *Journal of Clinical Oncology*, 12(10): 2086–2093.
- Selvin, S. 2008. *Survival analysis for epidemiologic and medical research*. Cambridge University Press.
- Sklar, A. 1959. Fonctions de répartition à n dimensions et leurs marges. *Publications de l'Institut Statistique de l'Université de Paris*, 8: 229–231.
- Tang, W.; Ma, J.; Mei, Q.; and Zhu, J. 2022. Soden: a scalable continuous-time survival model through ordinary differential equation networks. *Journal of Machine Learning Research*, 23(34): 1–29.
- Templeton, A. J.; Amir, E.; and Tannock, I. F. 2020. Informative censoring — a neglected cause of bias in oncology trials. *Nature Reviews Clinical Oncology*, 17(6): 327–328.
- Thomas Brooks, D. P. 1989. Airfoil self-noise.
- Tsiatis, A. 1975. A nonidentifiability aspect of the problem of competing risks. *Proceedings of the National Academy of Sciences*, 72(1): 20–22.
- V E, S.; Shin, C.; and Cho, Y. 2020. Efficient energy consumption prediction model for a data analytic-enabled industry building in a smart city. *Building Research & Information*, 49(1): 127–143.
- Voronov, S.; Frisk, E.; and Krysander, M. 2018. Data-driven battery lifetime prediction and confidence estimation for heavy-duty trucks. *IEEE Transactions on Reliability*, 67(2): 623–639.
- Wang, P.; Li, Y.; and Reddy, C. K. 2019. Machine learning for survival analysis. *ACM Computing Surveys*, 51(6): 1–36.
- Widder, D. V. 2010. *The laplace transform*. Dover Publications.

- Williams, J. S.; and Lagakos, S. W. 1977. Models for Censored Survival Analysis: Constant-Sum and Variable-Sum Model. *Biometrika*, 64(2): 215.
- Zhang, W.; Le, T. D.; Liu, L.; Zhou, Z.-H.; and Li, J. 2017. Mining heterogeneous causal effects for personalized cancer treatment. *Bioinformatics*, 33(15): 2372–2378.
- Zhang, W.; Li, J.; and Liu, L. 2021. A unified survey of treatment effect heterogeneity modelling and uplift modelling. *ACM Computing Surveys*, 54(8): 1–36.

A Survival Marginals

One of the primary goals in survival analysis is to estimate the survival function. The survival function S , defined as:

$$S(t) = Pr(T > t), \quad (10)$$

is a monotonically decreasing function characterizing the probability of subject survival past time t . A probability quantity closely related to the survival function is the distribution function $F(t) = 1 - S(t)$, which represents the probability of an event occurs before t .

Another desirable estimand of survival analysis is the hazard function $h(t)$. Intuitively, the hazard function is the instantaneous failure rate conditional on the subject has not failed until time t . Specifically, $h(t)$ is defined as

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{Pr(t \leq T < t + \Delta t | T \geq t)}{\Delta t} = \frac{f(t)}{S(t)}, \quad (11)$$

where $f(t) = \frac{\partial}{\partial t} F(t) = -\frac{\partial}{\partial t} S(t)$ denote the death density function. Furthermore, it is straightforward to represent the hazard function $h(t)$ as:

$$h(t) = \frac{f(t)}{S(t)} = -\frac{d}{dt} S(t) \cdot \frac{1}{S(t)} = -\frac{d}{dt} [\log S(t)]. \quad (12)$$

The cumulative hazard function is defined as the integral of the hazard function:

$$H(t) = \int_0^t h(u) du = -\log S(t). \quad (13)$$

Therefore, the survival function can be also expressed as $S(t) = \exp(-H(t))$.

B Censoring Mechanisms

We focus our discussion of censoring mechanisms on right-censoring. The censoring mechanism in most observational studies are unknown with few exceptions (Leung, Elashoff, and Afifi 1997). Specifically, there are two types of known censoring mechanism, namely Type I and Type II censoring. Type I censoring occurs in clinical observational studies when every subject is followed up for a specified time of C_0 or until event time. Type II censoring is often used in engineering experiments where a total of n devices are used in the study but, instead of continuing until all devices fail, the study is terminated when a total of pre-defined r devices fail. Generally speaking, Type II censoring does not affect estimation because it does not affect the likelihood model. However, Type II censoring is rarely adopted in applications except engineering. Both type I and type II censoring happen due to *end of study*.

However, as conducting trials and experiments are often expensive and time consuming, data with known censoring mechanism is rather difficult to collect. Therefore, researchers often have to be content with observational studies where the data is passively collected and censoring happens due to *lost of follow up*.

Therefore, in most applications that use observational data, the censoring mechanism are unknown and more complicated than the Type I and Type II designs. Within statistical literature, there are several widely used assumptions (Kleinbaum 2012) (we use intuitive definitions for readability, and refer readers to the referenced textbook for mathematical definitions) :

Definition 1. (Random Censoring) *Random censoring indicates that subjects who are censored at time t are representative of all subjects who remain at risk at time t with respect to the event of interest.*

Definition 2. (Independent Censoring) *Independent censoring indicates that within any subgroup of interest, subjects who are censored at time t are representative of all subjects in that subgroup who remain at risk at time t with respect to the event of interest.*

Definition 3. (Non-informative censoring) *Non-informative censoring indicates the distribution of survival time T provides no information about the distribution of censoring time U , and vice versa.*

Strictly speaking, independent censoring assumptions is a special case of non-informative censoring (Lagakos 1979). However, since examples where censoring is non-informative but not independent are artificially crafted (Williams and Lagakos 1977), we do not differentiate these two assumptions in our discussions.

Despite the importance of non-informative censoring assumption, it is impossible to identify whether the censoring mechanism is non-informative from observational data (Tsiatis 1975). Furthermore, as Kaplan and Meier (Kaplan and Meier 1958) noted, “in practice this assumption (independent censoring) deserves special scrutiny”. Therefore, it is necessary to relax the independence assumption and accommodate for dependent censoring.

C Derivation of Survival Likelihood

Without assumption on the censoring mechanism, we have

$$\begin{aligned}
\mathcal{L}_{\text{dep}} &= \Pr(T = t, U > t|X)^\delta \cdot \Pr(T > t, U = t|X)^{1-\delta} \\
&= \left\{ -\frac{\partial}{\partial y} \Pr(T > y, U > t)|_{y=t} \right\}^\delta \cdot \left\{ -\frac{\partial}{\partial z} \Pr(T > t, U > z)|_{z=t} \right\}^{1-\delta} \\
&= \left\{ \int_t^\infty f_{U|T,X}(u^*|t, X) f_{T|X}(t|X) du^* \right\}^\delta \cdot \left\{ \int_t^\infty f_{T|U,X}(t^*|u, X) f_{U|X}(u) dt^* \right\}^{1-\delta} \\
&= \left\{ f_{T|X}(t|X) \int_t^\infty f_{U|T,X}(u^*|t) du^* \right\}^\delta \cdot \left\{ f_{U|X}(t|X) \int_t^\infty f_{T|U,X}(t^*|u) dt^* \right\}^{1-\delta} \\
&= \left\{ f_{T|X}(t|X) \frac{\partial}{\partial u_1} \mathcal{C}(u_1, u_2) \Big|_{\substack{u_1 = S_{T|X}(t|X) \\ u_2 = S_{U|X}(t|X)}} \right\}^\delta
\end{aligned} \tag{14}$$

$$\cdot \left\{ f_{U|X}(t|X) \frac{\partial}{\partial u_2} \mathcal{C}(u_1, u_2) \Big|_{\substack{u_1 = S_{T|X}(t|X) \\ u_2 = S_{U|X}(t|X)}} \right\}^{1-\delta} \tag{15}$$

When conditional independent censoring is assumed, i.e., $T \perp\!\!\!\perp C|X$, the likelihood can be written as

$$\begin{aligned}
\mathcal{L}_{\text{indep}} &= \Pr(T = t, U > t|X)^\delta \cdot \Pr(T > t, U = t|X)^{1-\delta} \\
&= \{\Pr(T = t|X) \Pr(U > t|X)\}^\delta \cdot \{\Pr(T > t|X) \Pr(U = t|X)\}^{1-\delta} \\
&= \{f_{T|X}(t|X) S_{U|X}(t|X)\}^\delta \cdot \{f_{U|X}(t|X) S_{T|X}(t|X)\}^{1-\delta}
\end{aligned} \tag{16}$$

$$\propto f_T(t|X)^\delta S_T(t|X)^{1-\delta}. \tag{17}$$

Contrasting Equation 15 and 16, the potential bias caused by the independence censoring assumption is evident. The independent likelihood is only unbiased if the copula partial derivatives equal to the marginal survival functions, i.e., $\mathcal{C}(u_1, u_2) = u_1 u_2$. Since the dependency structure is usually unknown in observational studies and unverifiable from the data, the independent censoring assumption should be avoided in most practical applications.

D Copulas

Copulas are distribution functions of d -dimensional multivariate distributions with uniform margins. More formally, a function $\mathcal{C}(u_1, \dots, u_d) : [0, 1]^d \rightarrow [0, 1]$ is a copula if the following hold (Nelsen 1998):

- (Grounded) The copula equals to 0 if any of its argument is 0, i.e.,

$$\mathcal{C}(u_1, \dots, u_{i-1}, 0, u_{i+1}, \dots, u_d) = 0.$$

- (Uniform) The copula equals to u if one argument is u and all others are 1, i.e.,

$$\mathcal{C}(1, \dots, 1, u, 1, \dots, 1) = u.$$

- (d -increasing) For all $u = (u_1, \dots, u_d)$ and $v = (v_1, \dots, v_d)$ where $u_i < v_i$ for all $i \in [d]$,

$$\sum_{(w_1, \dots, w_d) \in \times_{i=1}^d \{u_i, v_i\}} (-1)^{|i:w_i=u_i|} \mathcal{C}(w_1, \dots, w_d) \geq 0$$

- (d -increasing) For each hyperrectangle $B = \Pi_{i=1}^d [u_i, v_i] \subseteq [0, 1]^d$, the C-volume of B is non-negative, i.e., $\int_B d\mathcal{C}(u) \geq 0$.

Intuitively, the d -increasing property states that the probability of any d -dimensional hyperrectangle is non-negative.

Copulas have been widely applied to areas such as quantitative finance and engineering thanks to *Sklar's theorem* (Sklar 1959), which states that any d -dimensional continuous joint distribution can be *uniquely* expressed with d uniform marginals and a copula \mathcal{C} . Formally, Sklar's theorem states that

Sklar's theorem (Sklar 1959) states that there exist a unique copula \mathcal{C} such that the joint distribution of event and censoring can be expressed as

$$F_{T,C}(t, c) = \mathcal{C}(F_T(t), F_C(c)). \tag{18}$$

According to probability integral transform, we know that $F_T(t)$ and $F_C(c)$ are uniform in $[0, 1]$. Note that instead of Equation 18 we can also assume a survival copula such that $P(T > t, C > c) = \tilde{\mathcal{C}}(1 - F_T(t), 1 - F_C(c))$ where $\tilde{\mathcal{C}}(u, v) = u + v - 1 + \mathcal{C}(1 - u, 1 - v)$. However, as there is no substantial difference between these two expressions (Nelsen 1998).

Algorithm 1: Survival Modeling with Unknown Dependent Censoring

Input: \mathcal{D} : right-censored survival dataset $\mathcal{D} = (X_i, t_i, \delta_i)$ for $i = 1, \dots, N$; \mathcal{M}_T : parametric survival model that emits the survival marginals $\hat{S}_{T|X}(t|X), \hat{f}_{T|X}(t|X)$; \mathcal{M}_U : parametric survival model that emits the censoring marginals $\hat{S}_{U|X}(t|X), \hat{f}_{U|X}(t|X)$; Learning rate α ; Number of iterations S
Output: $\varphi_{nn}, \psi_T, \psi_U$, the learned copula, survival, and censoring models

- 1: $\mathcal{M}_T \leftarrow$ Instantiate the survival marginal ($M_1; \psi_T$);
 - 2: $\mathcal{M}_U \leftarrow$ Instantiate the censoring marginal ($M_2; \psi_U$);
 - 3: **for** $i = 1, \dots, S$ **do**
 - 4: Compute \mathcal{L}_{dep} for a mini-batch of \mathcal{D} ;
 - 5: $\psi_T, \psi_U \leftarrow \text{AdamWUpdate}(\mathcal{L}_{\text{dep}}, (\psi_T, \psi_U), \alpha)$;
 - 6: $\Phi \leftarrow \text{AdamWUpdate}(\mathcal{L}_{\text{dep}}, \varphi_{nn}, \alpha)$;
-

E Theorem Proofs

For the sake of conciseness, we omit the covariate X from the notations. As we do not assume conditionally independent censoring, the derivations with covariates follow similarly.

Lemma 5. Suppose φ is differentiable on $(0, 1)$. If $\lim_{u \rightarrow 0} \varphi'(u) \in (-\infty, 0)$, then

$$\lim_{t \rightarrow 0} \frac{\partial}{\partial u_1} \mathcal{C}(u_1, u_2) \Big|_{\substack{u_1 = S_T(t) \\ u_2 = S_U(t)}} = 1.$$

Proof. From the definition of $S(t)$ we know that $\lim_{t \rightarrow 0} S_T(t) = \lim_{t \rightarrow 0} S_U(t) = 1$. Furthermore, we have $\lim_{t \rightarrow 1} \varphi^{-1}(t) = 0$ and $\varphi(0) = 1$ from the definition of the generator. Lastly, utilizing the definition of Archimedean copula, i.e., $\mathcal{C} = \varphi(\varphi^{-1}(u_1) + \varphi^{-1}(u_2))$, we can write

$$\begin{aligned} & \lim_{t \rightarrow 0} \frac{\partial}{\partial u_1} \mathcal{C}(u_1, u_2) \Big|_{\substack{u_1 = S_T(t) \\ u_2 = S_U(t)}} \\ &= \lim_{t \rightarrow 0} \frac{\varphi'(\varphi^{-1}(S_U(t)) + \varphi^{-1}(S_T(t)))}{\varphi'(\varphi^{-1}(S_T(t)))} = \lim_{t \rightarrow 0} \frac{\varphi'(t)}{\varphi'(t)} = \frac{c}{c} = 1 \end{aligned} \quad (19)$$

The case for $\partial \mathcal{C}(u_1, u_2) / \partial u_2$ can be similarly shown. □

Theorem 6. All copulas expressed by the neural network specified by φ^{nn} satisfy condition (C2) of Theorem 3.

Proof. The theorem can be proved by induction with respect to the indexes of the hidden layers. For notation simplicity, we assume that each hidden layer contains the same number of H hidden units. The results can be straightforwardly extended to the general cases. If φ^{nn} has one hidden layer, i.e. $L = 1$, the output from the last hidden layer of φ^{nn} can be straightforwardly expressed as

$$\varphi_1^{nn}(u) = \sum_{k=1}^H A_{1,j,k} e^{-B_{1,j} \cdot u}. \quad (20)$$

Since the output layer performs convex combination of $\varphi_1^{nn}(u)$, $B_{i,j}$ is positive, $A_{l,j,k}$ is non-negative with $\sum_j A_{l,i,j} = 1$, we can express $\varphi^{nn}(u)$ as $\varphi^{nn}(u) = \sum a \cdot \exp(-b \cdot t)$ with some $a, b > 0$. It follows naturally that $\varphi'(u) < 0$ for all u . Utilizing Lemma 4, the theorem is satisfied when $L = 1$.

Assume that $\varphi^{nn}(u)$ has L hidden layer, and the output of its last hidden layer can be expressed as $\varphi_L^{nn}(u) = \sum a \cdot \exp(-b \cdot t)$ with some $a, b > 0$, we have $\varphi'(u) < 0$ when $L = n + 1$ using the fact that the output layer is a convex combination of $\varphi_L^{nn}(u)$. Utilizing Lemma 4 again, the theorem is satisfied when for $L = n + 1$, completing the induction. □

F Additional Experiment Details and Results

Experiment Setting and Implementation Details

Experiments are conducted on a PC with one NVIDIA RTX4090 GPU. We utilize the Pytorch for implementing all neural networks and automatic differentiation. Tensors are computed with double precision (fp64) as the inversion of φ mandates numerical precision. When using Newton's method to compute the inverse φ_{nn}^{-1} , we terminate when the error is less than 1×10^{-12} .

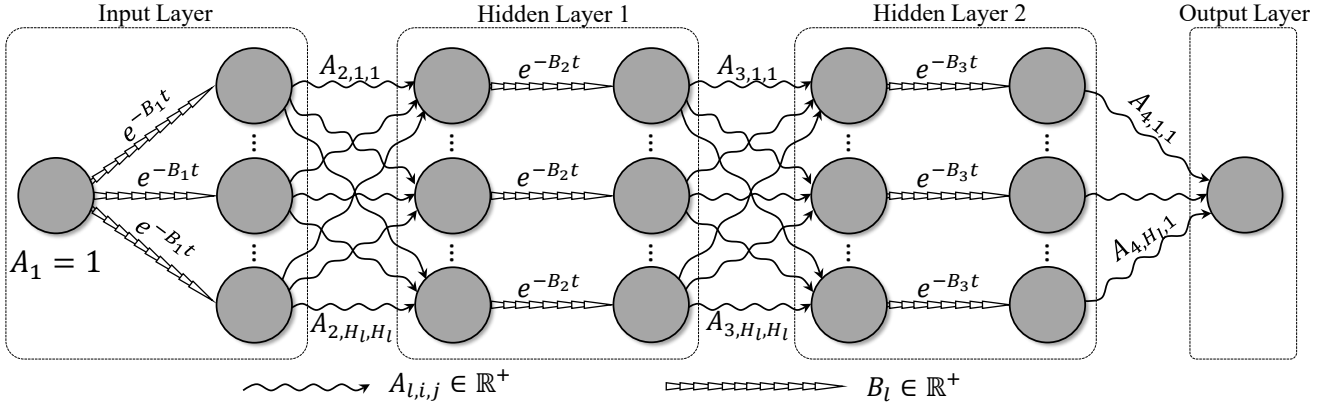


Figure 5: Network structure for representing the Archimedean copula generator φ_{nn} . The coefficients of A are parameterized by Softmax to ensure a convex combination. B is non-negative such that $\exp(-Bt)$ are negative exponentials.

Table 2: The Archimedean copulas used for generating dependence among survival and censoring outcomes and evaluating the performances of DCSURVIVAL. For Frank copula, D_1 is the Debye function of the first kind, i.e., $D_1(\theta) = \frac{1}{\theta} \int_0^\theta \frac{t}{e^t - 1} dt$.

Copula	$C_\theta(u, v)$	$\varphi_\theta(t)$	Kendall's τ
Clayton	$(u^{-\theta} + v^{-\theta} - 1)^{-\frac{1}{\theta}}$	$\frac{1}{\theta}(t^{-\theta} - 1)$	$\frac{\theta}{\theta+2}$
Gumbel	$\exp[-((\log(u))^\theta + (-\log(v))^\theta)^{\frac{1}{\theta}}]$	$(-\log t)^\theta$	$\frac{\theta-1}{\theta}$
Frank	$-\frac{1}{\theta} \log[1 + \frac{(e^{-\theta u} - 1)(e^{-\theta v} - 1)}{e^{-\theta} - 1}]$	$-\log \frac{e^{-\theta t} - 1}{e^{-\theta} - 1}$	$1 + \frac{4}{\theta}(D_1(\theta) - 1)$
Independent	uv	$-\log t$	0

For all our experiments we set φ_{nn} with $L = 2$ and $H_1 = H_2 = 10$, i.e., the copula representation contains two hidden layers with each of width 10. The network is small but sufficient for the dependency structure since the generator φ_{nn} is only 1-dimensional. Φ_B and Φ_B were uniformly initialized in the range $[0, 1]$ and $(0, 2)$. We use AdamW (Loshchilov and Hutter 2019) for optimization and use 30% of the training samples for evaluating the validation log-likelihood. No further hyperparameter tuning was performed. An illustration for the network of φ_{nn} can be found in Figure 5.

Algorithm 2: Synthetic Data Generating Process

Input: X : set of covariates; ν_T, ρ_T ; ν_U, ρ_U : parameters of the marginal distribution; $\psi_T(X), \psi_U(X)$: risk functions; C_θ : a pre-defined copula parameterized by θ .

Output: $\mathcal{D} = (X_i, t_i, \delta_i)$ for $i = 1, \dots, N$.

- 1: $\mathcal{D} = \emptyset$;
 - 2: **for** $i = 1, \dots, N$ **do**
 - 3: Sample $u_1 \sim \mathcal{U}_{[0,1]}$ and $u_2 \sim \mathcal{U}_{[0,1]}$ such that $F(u_1, u_2) = C_\theta(F(u_1), F(u_2))$;
 - 4: $T \leftarrow \left(\frac{-\log u_1}{\exp \psi_T(X)} \right)^{\frac{1}{\nu_T}} \cdot \rho_T$;
 - 5: $U \leftarrow \left(\frac{-\log u_2}{\exp \psi_U(X)} \right)^{\frac{1}{\nu_U}} \cdot \rho_U$;
 - 6: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(X_i, \min(T_i, U_i), \mathbb{I}_{[T_i < U_i]})\}$.
-

We implement DCSURVIVAL with PyTorch 2.0.0. The pseudo code for the algorithm procedure can be found in Algorithm 1. DCSURVIVAL uses the same hyper-parameters across all experiments except the learning rate. We set the learning rate to $1e - 4$ for the synthetic datasets and the semi-synthetic datasets STEEL and Airfoil. For the real-world dataset GBSG2 and SEER, we use the learning rate of $1e - 5$. All experiments are conducted with an NVIDIA RTX4090 GPU where the average training speed is approximately 3 epoch/s. However, as DCSURVIVAL uses double precision tensors, access to professional-grade GPU such as a NVIDIA V100 (7.8 TFLOPS for V100 vs 1.3 TFLOPS for RTX4090) will significantly boost the training speed.

For implementation of the compared algorithms, we use the PYCOX package (Kvamme, Ørnulf Borgan, and Scheel 2019)

Algorithm 3: Semi-Synthetic Censoring Inducing Procedure

Input: X : set of covariates; T : set of dependent variables; \mathcal{C} : A copula specifying the dependency structure.

Output: $\mathcal{D} = (X_i, t_i, \delta_i)$ for $i = 1, \dots, N$. A censored dataset where the joint distribution of T and U are governed by the copula \mathcal{C} .

```
1:  $M_T \leftarrow \text{Weibull}(X, Y, \mathbf{1}^N)$ ; Fit event model using Weibull distribution
2:  $M_U \leftarrow M_T$ 
3:  $M_U.\nu \leftarrow M_T.\nu/0.8$ ; Reduce variance for the censoring distribution
4:  $\mathcal{D} = \emptyset$ ;
5: for  $i = 1, \dots, N$  do
6:    $u_i \leftarrow S_T(Y_i)$ ; obtain event quantile
7:    $v_i \sim \mathcal{C}(\cdot | u_i)$  Conditionally sample censoring quantile from  $\mathcal{C}$ 
8:    $U_i \leftarrow S_U(v_i)$  Obtain censoring time with probability inverse transform
9:    $\mathcal{D} \leftarrow \mathcal{D} \cup \{(X_i, \min(T_i, U_i), \mathbb{I}_{[T_i < U_i]})\}$ .
```

which implements DEEPSURV, DEEPHIT, COXCC, COXTIME and PCHAZARD. For RANDOM SURVIVAL FOREST and COXPH methods, we utilize the implementation provided in SCIKIT-SURVIVAL (Pölsterl 2020). The code of KNOWNCOPULA can be found at https://github.com/rgklab/copula_based_deep_survival. For classical survival methods such as RANDOM SURVIVAL FOREST and COXPH, we use the default parameters. For all methods that utilize deep neural networks, we use the same multi-layer perceptron (MLP) structure of [32, 32, 32] and *ReLU* activation function. The same network structure is also used for the covariate representation in DCSURVIVAL, albeit DCSURVIVAL uses the hyperbolic tangent activation function. A validation set is used for all deep learning based algorithms, which consists of 30% of the training samples.

To ensure a fair comparison among all algorithms, we use the random seed for all the data generation procedure and train/val/test set splits. In other words, for each of the ten repeats, every algorithm use the same training, validation and test samples.

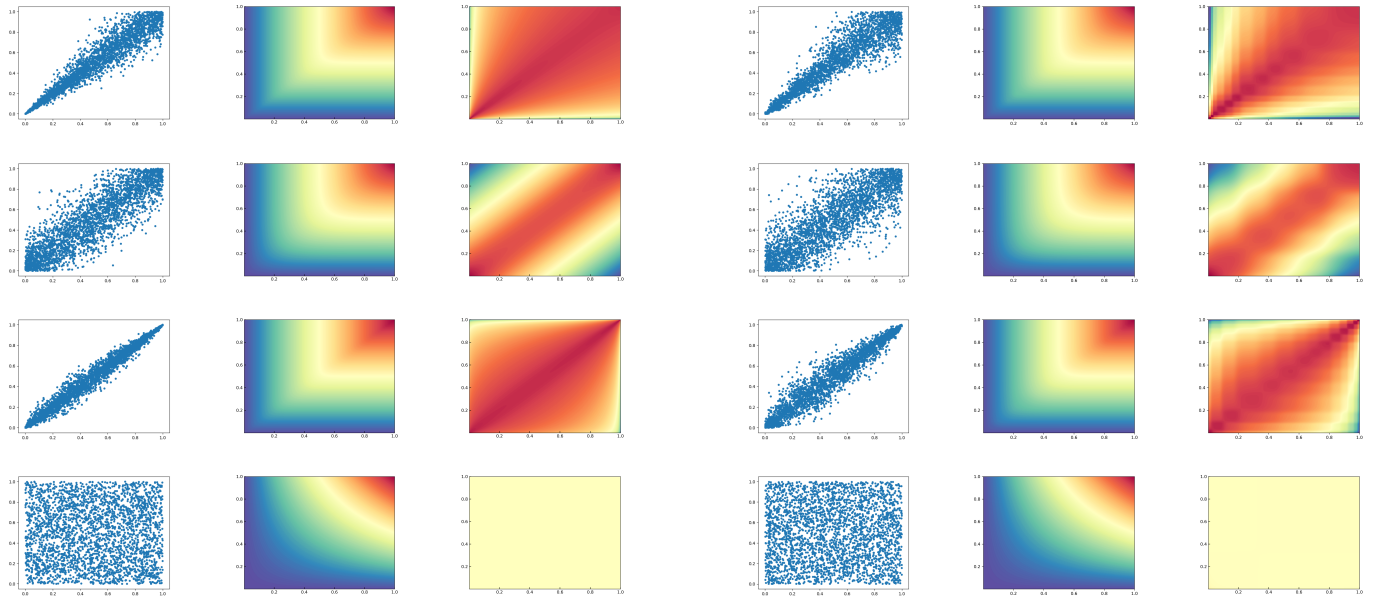
Datasets

We provide the pseudocode for generating the synthetic datasets in Algorithm 2, and the pseudocode for inducing censoring in the semi-synthetic datasets in Algorithm 3. Furthermore, the details of different Archimedean copulas used during simulation and the relationships between θ and Kendall's τ are provided in Table 2.

For the real-world GBSG2 dataset, it can be accessed using the LIFELINES package (Davidson-Pilon 2023). The SEER dataset can be accessed using the SEERSTAT software downloaded from <https://seer.cancer.gov/seerstat/>.

Results from Learning Dependence under Non-linear Risk

We present the scatter plot, cumulative distribution plot, and log-probability density plots for learning the dependency structure from right-censored observations with non-linear risks in Figure 6. It can be seen that DCSURVIVAL successfully learns the underlying copulas governing the dependency, as the samples from learned models closely resemble the ground truth ones.



(a) Ground Truth

(b) Learned Copula

Figure 6: Top to bottom: Learning Clayton, Frank, Gumbel, and Independence copulas using DCSURVIVAL from right-censored observations with non-linear survival and censoring risks. Plots from left to right: (i) samples drawn from the ground truth and learned distributions. (ii) joint cumulative distributions, (iii) log probability densities.