# SurvivalGAN: Generating Time-to-Event Data for Survival Analysis

**Alexander Norcliffe**[*]
University of Cambridge
alin2@cam.ac.uk

**Bogdan Cebere**[*]
University of Cambridge
bcc38@cam.ac.uk

**Fergus Imrie**
University of California, Los Angeles
imrie@ucla.edu

**Pietro Liò**
University of Cambridge
pl219@cam.ac.uk

**Mihaela van der Schaar**
University of Cambridge
Alan Turing Institute
mv472@cam.ac.uk

## Abstract

Synthetic data is becoming an increasingly promising technology, and successful applications can improve privacy, fairness, and data democratization. While there are many methods for generating synthetic tabular data, the task remains non-trivial and unexplored for specific scenarios. One such scenario is survival data. Here, the key difficulty is censoring: for some instances, we are not aware of the time of event, or if one even occurred. Imbalances in censoring and time horizons cause generative models to experience three new failure modes specific to survival analysis: (1) generating too few at-risk members; (2) generating too many at-risk members; and (3) censoring too early. We formalize these failure modes and provide three new generative metrics to quantify them. Following this, we propose SurvivalGAN, a generative model that handles survival data firstly by addressing the imbalance in the censoring and event horizons, and secondly by using a dedicated mechanism for approximating time-to-event/censoring. We evaluate this method via extensive experiments on medical datasets. SurvivalGAN outperforms multiple baselines at generating survival data, and in particular addresses the failure modes as measured by the new metrics, in addition to improving downstream performance of survival models trained on the synthetic data.

## 1 INTRODUCTION

Deep learning has seen incredible success in recent years, yet generally deep models still require large amounts of high-quality data to train well. Data collection is expensive, and often privacy constraints limit how much data can actually be used or shared (De Capitani Di Vimercati et al., 2012; Jain et al., 2016). Synthetic data generation is a popular solution that aims to create new data that mirrors the statistical properties of the original dataset, tackling the need for privacy and more data in parallel (Jordon et al., 2022). Synthetic data has significant promise, with the potential to improve: (1) fairness & bias by generating data from under-represented groups (van Breugel et al., 2021); (2) robustness by augmenting an original dataset (Perez and Wang, 2017); (3) privacy by not using identifiable data to train a supervised model (Zhang et al., 2017; Jordon et al., 2018; Yoon et al., 2020), and (4) data democracy by allowing researchers with fewer resources to access inexpensive data (Benedetti et al., 2020; Wang et al., 2021). As a result, significant attention has been placed on developing generative models (Prakash et al., 2019; Tobin et al., 2017). Prominent examples can be found across many domains such as images (Karras et al., 2019; Hinterstoisser et al., 2019), audio (Oord et al., 2016), and medicine (Chen et al., 2021; Singh and Mukhopadhyay, 2011). One area that remains vastly unexplored is survival data.

Survival analysis seeks to answer the question: given some measurements at a fixed point in time, how long will it take until a specific event occurs? In engineering for example, given a machine's current condition when do we expect there to be mechanical failure (de Cos Juez et al., 2010), or in finance if a company's stock price is at a certain value how long will it be before they declare bankruptcy. Survival models are incredibly impactful; for instance in medicine (Lee and Go, 1997; Arsene and Lisboa, 2007) they can be used to estimate how long a patient is expected to survive

with a given disease such as Covid-19 (Kaso et al., 2022; Lu et al., 2021; Salinas-Escudero et al., 2020; Ali et al., 2022). They can also be used in clinical trials (Singh and Mukhopadhyay, 2011) to investigate how long it may be before a death, relapse, or adverse reaction. Further application areas include economics (Danacica and Babucea, 2010; LeClere, 2005) and sociology (Kent, 2010; Gross et al., 2014). Given these models' impact, we wish to generate the highest quality synthetic data for survival analysis.

Naively, generating synthetic data for survival analysis seems straightforward. However, there are two significant obstacles: tabular data and censored data. The tabular setup is notably more complex than image or text data where generative models have typically been applied (Xu et al., 2019). This is due to the mixture of categorical and continuous features and their joint distributions. On top of tabular complexity, in most cases survival data is not completely observed: we may not know when an event occurred – if at all. This data is said to be "censored". For instance in a drug trial, if the event is being cured of a disease and a subject withdraws from the trial, we will not know when, or even if, they were cured. As a result, the subject is censored at the point of withdrawal. In spite of missing the time of the event (if indeed an event occurred), censored data still contains information, we know that the event has *not occurred* before the time of censoring. Further, censored data is often more abundant than non-censored, hence, a good generative model will incorporate this data and the imbalance, despite the missing event.

**Contributions.** In this work we present a synthetic generation framework for efficiently handling censored tabular survival data. Our contributions are three-fold:

1. We formalize the synthetic data generation problem for survival analysis, identify three failure modes unique to the survival setting, and introduce three metrics to quantitatively evaluate these failures and provide a clearer understanding of the synthetic data's utility (Sections 3 & 4).
2. We propose SurvivalGAN, a method that is able to efficiently incorporate censored data, tabular data, and censoring & time-horizon imbalance to generate synthetic data to train survival models (Section 5).
3. We investigate SurvivalGAN via extensive experiments on five medical datasets. We demonstrate its ability to generate high-quality synthetic data relative to five robust benchmarks, as measured by established generative metrics, the new survival specific metrics, and downstream model performance (Section 6).

## 2 RELATED WORK

**Generative Models.** Generative models come in various flavors. Classically, Bayesian networks represent a high di-

mensional distribution with a directed acyclic graph to compactly give the dependency structure (Niedermayer, 2008). These can often be slow to sample from, requiring methods such as Markov Chain Monte-Carlo. More recent generative models include Variational Auto-Encoders (Kingma and Welling, 2013; Vahdat and Kautz, 2020), Generative Adversarial Networks (Goodfellow et al., 2014) and Normalizing Flows (Rezende and Mohamed, 2015; Kingma and Dhariwal, 2018). Exact training techniques and architectures differ, but typically modern models use deep networks to learn a mapping from an easy-to-sample latent space (such as Gaussian noise) to the data distribution, leading to fast sampling. The exception is diffusion-based models (Song and Ermon, 2019; Ho et al., 2020; Ramesh et al., 2022) which learn a reverse diffusion process in observation space to take noisy points to regions of higher probability via a series of (possibly many) function evaluations.

In their standard implementations, these methods are not well suited to tabular data, due to the mixture of continuous and categorical variables. Variants have been proposed that are designed to work in the tabular domain. Tabular GAN (Xu and Veeramachaneni, 2018) works on tabular data and CTGAN (Xu et al., 2019) is a GAN designed to work specifically with tabular data that mixes categorical and continuous variables. This is achieved with a tabular encoding and training by sampling. TabFairGAN (Rajabi and Garibay, 2022) extends this to include a fairness constraint generating accurate and fair data. Tabular variants of other models also exist, such as TVAE (Xu et al., 2019), RTVAE (Akrami et al., 2022), and GOGGLE (Liu et al., 2023), which adapt VAEs, while Vahdat et al. (2021) adapt score-based models to generate tabular data by running the diffusion in latent space. These principled methods are well suited to tabular data, however, survival data contains the added difficulty of censored data which these models are not able to handle.

**Survival Data.** Few models exist to generate synthetic survival data. Existing survival models learn to sample from the conditional distribution of event time given the initial state (known as the covariates), for example Bender et al. (2005) and Austin (2012) present statistical models that transform samples from a uniform distribution to survival times by inverting the cumulative hazard function, conditioned on the covariates. Sampling from $p(t|\boldsymbol{x})$ is a probabilistic survival model and does not generate the covariates, which we desire to make a full synthetic dataset. More modern techniques incorporate deep learning in the generative process. Ranganath et al. (2016) use deep exponential families to generate survival data, limiting the flexibility of the learnt distribution. Miscouridou et al. (2018) and Zhou et al. (2022) relax this assumption but still generate survival times and censoring statuses conditioned on covariates, rather than generating covariates and times. None of these models are able to generate survival data that considers censoring and

still generates covariates with high fidelity.

# 3 PROBLEM FORMULATION AND BACKGROUND

## 3.1 Generating Synthetic Survival Data

Synthetic data generation is concerned with creating new data points that resemble but are not identical to real data. Given a dataset of $N$ training observations, $\mathcal{D} = \{\boldsymbol{z}_i\}_{i=1}^N$ drawn from an unknown distribution $\boldsymbol{z}_i \sim p_{\boldsymbol{z}}$, the task is to generate $M$ new observations $\mathcal{D}^{\text{Syn}} = \{\boldsymbol{z}_j^{\text{Syn}}\}_{j=1}^M$ that appear to have been sampled from the same distribution. Instances from survival data are of the form $\boldsymbol{z}_i = (\boldsymbol{x}_i, t_i, E_i)$. Here, $\boldsymbol{x} \in \mathcal{X}$ are $m$-dimensional tabular covariates giving a subject's state at an initial time, containing continuous and categorical features. $t_i \in \mathcal{T}$ is the time of a given event, usually the initial time is 0 such that $\mathcal{T} = \mathbb{R}_{\geq 0}$. Finally, $E \in \mathcal{E}$ is the event indicator, typically $\mathcal{E} = \{0, 1\}$, where $E = 0$ means the individual is censored at $t$, and $E = 1$ means the event of interest occurred at $t$. In this exposition we consider a single event, but this can be naturally extended to competing events by extending $\mathcal{E}$ to include new discrete events.

## 3.2 Survival Analysis

Here we give a brief summary of classical survival analysis[1] relevant to this work. This is meant as a way to introduce key terms from survival analysis that we will use throughout the paper; for a thorough introduction see Jenkins (2005) and Machin et al. (2006).

**Survival function.** We can conduct survival analysis with the probability density function $p(t|\boldsymbol{x})$, giving the likelihood that the event of interest happens at time $t$ given covariate $\boldsymbol{x}$. With this, the **survival function** is defined as the probability that the event *has not* occurred by time $t$,

$$S(t|\boldsymbol{x}) = \int_t^\infty p(t'|\boldsymbol{x})dt'. \tag{1}$$

Equation (1) gives us the proportion of the subjects with covariate $\boldsymbol{x}$ that have survived up to point $t$. When the initial time is zero, the event cannot happen before $t = 0$, so $S(0|\boldsymbol{x}) = 1$, and $p(t|\boldsymbol{x})$ is a valid probability distribution (non-negative) so $S(t|\boldsymbol{x})$ is a decreasing function.

**Time-to-event approximation.** This is the task of approximating the expected lifetime for any $\boldsymbol{x}$. By definition this is given by $\mu(\boldsymbol{x}) = \int_0^\infty t'p(t'|\boldsymbol{x})dt'$. Via integration by parts, this is equal to the area under the survival curve $\mu(\boldsymbol{x}) = \int_0^\infty S(t|\boldsymbol{x})dt$. The above properties can be calculated at a population level by marginalizing out $\boldsymbol{x}$. Survival

---

[1] Predicting the time of event given a covariate, not generating survival data itself.

models usually fall into one of these two categories, i.e. (1) estimate the survival function, or (2) estimate the time-to-event. The first category is well-studied, with solutions ranging from linear models (Cox, 1972) to random forests (Ishwaran et al., 2008) to gradient boosting (Barnwal et al., 2022). The second category can be derived from the survival probabilities, but has also been studied independently using neural networks (Lee et al., 2018; Chapfuwa et al., 2018). Finally, it is possible to combine approaches from both categories via ensembling methods (Imrie et al., 2022).
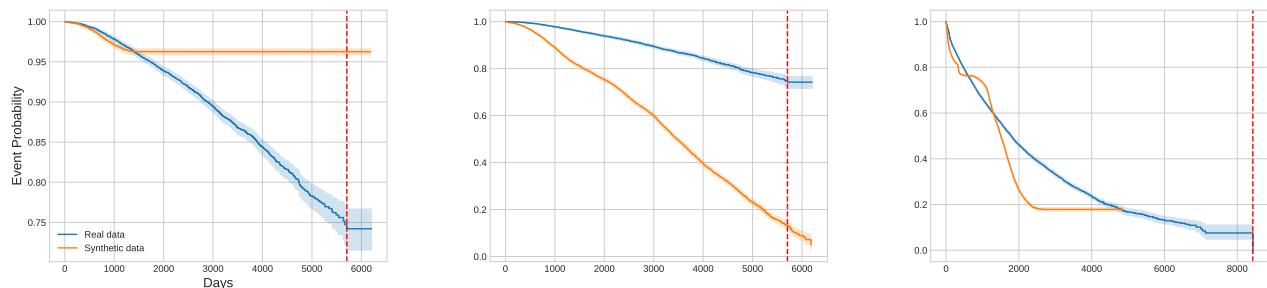
## 3.3 Challenges of Generating Survival Data

**Tabular Data.** Handling tabular data is non-trivial for most generative models. Some of the reasons are the mixed feature types (categorical or continuous), non-Gaussian distribution of the features, or highly imbalanced categorical features (Xu et al., 2019).

**Time and Censoring Failure Modes.** We are attempting to generate data from the full distribution $p(\boldsymbol{x}, t, E)$; generating covariates $\boldsymbol{x}$ is subject to the traditional failure modes of generative modeling. However, modeling the marginal distribution of the time and event pair $p(t, E)$ is specific to survival analysis and introduces three new failure modes, driven by two key dataset imbalances:

1. **Censoring imbalance:** There is often an imbalance in the amount of censoring. By not incorporating this balance, a model could generate unusable data, for example consisting of too many censored records. This makes a model **over-optimistic** because the event in question occurs less frequently than expected and the model predicts too high an expected lifetime; this is the first failure mode. On the other hand, by not generating enough censored data, the predicted lifetime could become too low, resulting in the model being **over-pessimistic**; this is the second failure mode.

2. **Time imbalance:** It is also possible for datasets to have a time-horizons imbalance. Where the majority of events are early, $t$ is early regardless of $E$. A model is at risk of learning this imbalance and making events too early. We don't want to focus only on short-term examples we want to have a broad view of the timeline. Such a model is **short-sighted**; this is the third failure mode.

We illustrate these failure modes in Figure 1. These show the Kaplan-Meier (KM) curves (Kaplan and Meier, 1958), which plot an approximation of the survival function at a given time. For a finite sample size, the KM curves show the proportion of the subjects that have made it to time $t$ without the event happening. We assume the sample size is large enough, such that the curve approaches the true survival function. To quantify these failure modes we introduce three metrics in the next section.

(a) Over-optimism. The area under the synthetic curve is larger than the true curve, the model is over-optimistic.

(b) Over-pessimism. The area under the synthetic curve is lower than the true curve, the model is over-pessimistic.

(c) Short-sightedness. The synthetic curve ends noticeably earlier than the true curve, the model is short-sighted.

Figure 1: The three possible failure modes specific to survival data, illustrated with Kaplan-Meier plots.

# 4 ASSESSING SYNTHETIC SURVIVAL DATA

Various metrics exist to evaluate the similarity between synthetic data and real data. For example, Maximum Mean Discrepancy (Sutherland et al., 2017), Inception Score (Salimans et al., 2016), and Fréchet Inception Distance (Heusel et al., 2017) are used to evaluate sample quality. However, representing the performance of a synthetic dataset with one metric is prone to over-simplifying the situation. In Alaa et al. (2022), three further model-independent metrics are proposed to overcome this, $\alpha$-Precision $\beta$-Recall and Authenticity. These are designed to evaluate sample quality, diversity, and similarity to real data, respectively. Importantly, none of these metrics are designed for survival data and are unable to capture the failure modes that originate from censoring imbalance and time imbalance.

We introduce three new metrics specific to evaluating synthetic survival data, targeting the failure modes in modeling the marginal distribution of $p(t, E)$. We have seen that as well as traditional generative failure modes (e.g. mode collapse or low sample quality), survival data presents three new ones: over-optimism, over-pessimism, and short-sightedness. We quantify these phenomena by looking at the differences in survival probabilities. Survival curves have been used to quantify the performance of predictive survival models but have not been used to evaluate synthetic survival data.

**Optimism.** The total area under the Kaplan-Meier plot gives us the expected lifetime of the population $\mu = \int_0^\infty S(t)dt$. We call the model over-optimistic if $\mu_{\text{Syn}} > \mu_{\text{Real}}$, and over-pessimistic for the reverse. Therefore we are interested in the quantity $\mu_{\text{Syn}} - \mu_{\text{Real}}$. The KM plots are only available in a finite interval, after which censoring is present. This can make the expected lifetime diverge if the survival function has not reached zero by this point. Instead, we consider the plots up to the latest available time

$T$. If the synthetic KM plot ends before the real one (due to short-sightedness), it is extrapolated assuming a constant rate of events. This gives the definition of optimism

$$\text{Optimism} = \frac{1}{T} \int_0^T \big( S_{\text{Syn}}(t) - S_{\text{Real}}(t) \big) dt. \quad (2)$$

This is the mean difference between the two Kaplan-Meier plots, and can also be viewed as a scaled difference in their areas. This metric takes values between -1 and 1, with 0 predicting the exact same expected lifetime, not over-optimistic or over-pessimistic. Positive values represent a synthetic expected lifetime higher than the true one, making the data over-optimistic and vice versa.

**Short-Sightedness.** Models trained on synthetic data may not be able to predict past a certain time horizon, meaning that the synthetic data is censored from that point. To quantify this using the KM plots, we consider the two end times $T_{\text{Syn}}$ and $T_{\text{Real}}$ and take their relative difference:

$$\text{Short-Sightedness} = \frac{T_{\text{Real}} - T_{\text{Syn}}}{T_{\text{Real}}}. \quad (3)$$

This quantifies the relative amount that the time horizons in the generated data stop before the real data. This metric takes values between 0 and 1, with 0 giving no short-sightedness and 1 giving full short-sightedness.

This metric can be generalized to also measure long-sightedness, where the predicted times are larger than the true times. Instead of dividing by $T_{\text{Real}}$, we divide by $\max(T_{\text{Real}}, T_{\text{Syn}})$, giving a value between -1 and 1, 0 being perfect, 1 being maximally short-sighted and -1 being maximally long-sighted. We observed that this does not happen in practice and therefore use the simpler definition in Equation (3).

**Kaplan-Meier Divergence.** It is possible to have scores of zero for both optimism and short-sightedness but still have non-matching Kaplan-Meier curves. Therefore we

finally include the mean absolute difference between the curves, which we call the Kaplan-Meier divergence (KM Divergence),

$$\text{KM Divergence} = \frac{1}{T} \int_0^T \big| S_{\text{Syn}}(t) - S_{\text{Real}}(t) \big| dt. \quad (4)$$

This will be between 0 and 1 (because $S(t)$ is always between 0 and 1), with 0 when the curves match perfectly, and 1 when they have the maximum difference possible at all times. In Appendix B, we demonstrate the need for all three metrics, showing they target different ways the KM curves can differ; with optimism and short-sightedness having an interpretable meaning. We also show that it is possible to bound the optimism by the total-variation divergence of the underlying probability density functions.

**Other Metrics.** Besides measuring the quality of $p(t, E)$ using these new metrics, we can measure the quality of the covariate marginal distribution $p(\boldsymbol{x})$ using a standard generative metric of choice. To measure the quality of the full distribution $p(\boldsymbol{x}, t, E)$, we evaluate the downstream performance of models trained with synthetic data compared to those trained with real data.

## 5 OUR MODEL - SURVIVALGAN

Below we describe our solution – SurvivalGAN – in depth. A block diagram is given in Figure 2. Briefly, to generate synthetic data, a condition vector $C$ and event $E$ are given by the user. The condition is a one-hot vector with both interpretable features (sex for example) and latent encodings (see full description below) and the event indicates censoring or true event (0 or 1). These can either be sampled according to the training data frequencies ($C, E \sim p_{C,E}$) or by manually selecting them. A conditional GAN is used to generate a covariate $\boldsymbol{x} \sim p_{\boldsymbol{x}|C,E}$. Finally, a survival function and time-to-event/censoring regressor are used together to generate the time $t \sim p_{t|\boldsymbol{x},E}$. One key insight is that the method assumes that censoring occurs independently and at random, and given some new covariates and a censoring status, a separate model determines the time of event/censoring. This allows us to follow a similar censoring ratio as the training dataset, without risking censoring all the synthetic instances.

### 5.1 Model and Training

**Conditional GAN.** The conditional GAN (part 1 of Figure 2), allows us to sample $\boldsymbol{x} \sim p_{\boldsymbol{x}|C,E}$, where $C$ is a user determined condition that the generator takes as input. We start by training a tabular encoder. This is critical as: (1) it allows us to handle censoring and time imbalance (as well as imbalance in the covariates) during training and generation and (2) it enables us to handle continuous and categorical variables; these were the key difficulties specific to survival data as described in Section 3.3. We follow the approach by

Xu et al. (2019). For each continuous feature, a Gaussian Mixture Model (GMM) (Reynolds, 2009) with $N_C$ components is trained. The tabular encoder for that feature is then given by Encoder : $\mathbb{R} \rightarrow \{0, 1\}^{N_C} \times \mathbb{R}$, where the first part of the output is a one-hot encoding telling us which component of the GMM the feature belongs to, and the second part is the number of standard deviations away from that component's mean $(x - \mu)/\sigma$. Note that the inverse of this encoding is trivial and does not need to be learnt since we include the one hot encoding of which mode the feature is in, and we know $\mu$ and $\sigma$ for that mode. For categorical features, the encoding is simply a one-hot encoding of the feature. Following this we also define a class encoder, ClassEncoder : $\mathbb{R} \rightarrow \{0, 1\}^{N_C}$, which simply takes the one-hot vector from the tabular encoding, saying which GMM component the feature belongs to without the location within that component. This allows us to represent a condition, $C$, for the generator. For extensive information on the tabular encoding see Reynolds (2009).

Once the tabular encoder has been trained, we train the GAN, which consists of generator $G_{\theta_g}$ and discriminator $D_{\theta_d}$. For a given sample $(\boldsymbol{x}, t, E)$ from the training dataset, the covariate is encoded first, $\boldsymbol{x}_e = \text{Encoder}(\boldsymbol{x})$. Following this, the condition is encoded from the input $C = \text{ClassEncoder}(\boldsymbol{x}, t, E)$. We then sample $\epsilon \sim \mathcal{U}_{[0,1]}$ and $\boldsymbol{z} \sim \mathcal{N}(0, I)$ and train with the Wasserstein GAN losses with gradient penalty:

$$
\begin{aligned}
\tilde{\boldsymbol{x}}_e &= \epsilon \boldsymbol{x}_e + (1 - \epsilon) G_{\theta_g}(C, \boldsymbol{z}) \\
L_G &= -D_{\theta_d}(C, G_{\theta_g}(C, \boldsymbol{z})) \\
L_D &= D_{\theta_d}(C, G_{\theta_g}(C, \boldsymbol{z})) - D_{\theta_d}(C, \boldsymbol{x}_e) \\
&\quad + \lambda(\|\nabla_{\tilde{\boldsymbol{x}}_e} D_{\theta_d}(C, \tilde{\boldsymbol{x}}_e)\|_2 - 1)^2,
\end{aligned}
$$

where $\lambda$ is the gradient penalty. We found that using the Wasserstein GAN with gradient penalty and the tabular encoder made training more stable.

**Survival Function.** The survival function (part 2 of Figure 2), is used with a time regressor (see next model component) to sample from $p_{t|\boldsymbol{x},E}$. The survival function $S : \mathcal{X} \times \mathcal{T} \rightarrow [0, 1]$ is a decreasing function predicting the survival probability at a given time horizon for given covariates $\boldsymbol{x}$. This is used in the generation process to create a set of survival probabilities at $N_H$ time horizons $\{S(\boldsymbol{x}, t_i)\}_{i=1}^{N_H}$. In practice, we use 100 evenly spread times between the minimum and maximum horizons in the training data. Any survival model can be used in this setup, thus the training is general. We use DeepHit (Lee et al., 2018) because, beyond strong predictive performance, it can be extended to competing risks, making it a flexible solution for future extensions. DeepHit uses a custom loss function consisting of two terms, one for log-likelihood of the joint distribution of $t$ and $E$, and one incorporating cause-specific losses. We refer the reader to Lee et al. (2018) for more detail.
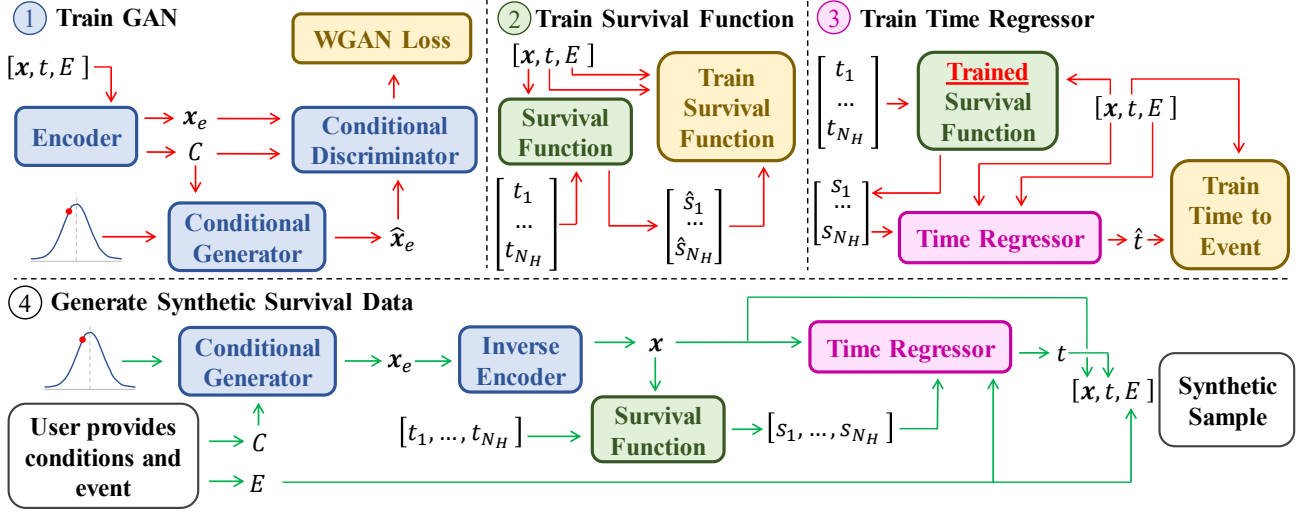
Figure 2: Block diagram of SurvivalGAN's components, their interaction, and the generation process. Both the specific survival function and time regressor may be chosen, hence the generic training for those components. Arrows pointing to the parentheses around vectors imply the whole vector is the input/output of a function. Arrows pointing to specific vector components only use those.

**Time Regressor.** The time regressor (part 3 of Figure 2) is trained after the survival function as it relies on using a trained survival model. It is used with the survival function to sample from $p_{t|\boldsymbol{x},E}$. The time regressor is a function $T : \mathcal{X} \times [0,1]^{N_H} \times \mathcal{E} \to \mathcal{T}$, which takes as input a covariate, event type and $N_H$ outputs of the survival function. The function then predicts the time that the event or censoring happens such that $T(\boldsymbol{x}, S(\boldsymbol{x}, t_1), S(\boldsymbol{x}, t_2), ..., S(\boldsymbol{x}, t_{N_H}), 0)$ gives the time to *censoring* for a given covariate, and $T(\boldsymbol{x}, S(\boldsymbol{x}, t_1), S(\boldsymbol{x}, t_2), ..., S(\boldsymbol{x}, t_{N_H}), 1)$ is the time to *event* for a given covariate. Any time regressor may be used, hence the training is general. We use XGBoost (Chen and Guestrin, 2016), with the mean squared error of $\log(t)$ as the loss.

### 5.2 Generation

Generation from the model is shown in part 4 of Figure 2. First, a one-hot condition vector $C$ and type of event $E$ are provided by the user. The user can choose these, allowing us to sample underrepresented groups within the covariates, and target specific edge cases when creating a synthetic dataset if desired. Alternatively, and the method we opt for in our empirical evaluation, one can sample conditions from a categorical distribution constructed from training data frequencies, which we call the **Imbalanced Sampler**. The generator uses this condition vector to produce an encoded covariate $\boldsymbol{x}_e = G_{\theta_g}(C, \boldsymbol{z})$, where $\boldsymbol{z} \sim \mathcal{N}(0, I)$. This is deterministically inverted to produce the covariate $\boldsymbol{x} = \text{Encoder}^{-1}(\boldsymbol{x}_e)$. Finally, this is used in the survival function with predefined time horizons and the time regressor to generate the time of event/censoring $t = T(\boldsymbol{x}, S(\boldsymbol{x}, t_1), S(\boldsymbol{x}, t_2), ..., S(\boldsymbol{x}, t_{N_H}), E)$.

**Necessity of Components.** To demonstrate that all components of SurvivalGAN are necessary, we carry out an ablation study in Section 6.3 individually testing each module. Table 3 demonstrates that all are crucial for the quality of the synthetic survival data.

## 6 EXPERIMENTS

To assess the quality of SurvivalGAN, we evaluate the following aspects:

1. **Quality of Marginals:** Section 6.1 analyses how closely the distribution of synthetic samples matches the original from two perspectives: (1) the marginal distribution of covariates $p(\boldsymbol{x})$. This is done using Jensen-Shannon distance and Wasserstein distance between real and synthetic covariates. We do not necessarily expect SurvivalGAN to perform better than the baselines here, but we check it does not perform worse. (2) the censoring and temporal marginal $p(t, E)$, evaluated using the optimism, KM divergence, and short-sightedness metrics. We examine t-SNE plots (van der Maaten and Hinton, 2008) to provide a qualitative performance of the covariates and Kaplan-Meier plots for time and censoring variables in Appendix E.

2. **Downstream Performance:** Section 6.2 compares the performance of survival models trained with synthetic data generated by SurvivalGAN to models trained with data generated from the baselines. This is used to quantify the quality of the full $p(\boldsymbol{x}, t, E)$ distribution. Here, a good result is when a model trained with synthetic data performs similarly to one trained with real data or in rare cases better (Luo and Lu, 2018), and outperforms models trained with different synthetic data.

Table 1: The mean and standard deviation of the covariate and time-censoring metrics. The best values are in **bold**, for short-sightedness extreme failure cases are underlined. ∗: Not evaluated.

| Metric | Method | AIDS | CUTRACT | PHEART | SEER | METABRIC |
|---|---|---|---|---|---|---|
| Jensen-Shannon Distance (Lower is Better) | SurvivalGAN | **0.012**±**0.02** | 0.024±0.01 | 0.013±0.01 | 0.022±0.01 | 0.015±0.01 |
| | PrivBayes | 0.028±0.01 | 0.020±0.01 | **0.008**±**0.01** | 0.022±0.01 | 0.043±0.01 |
| | ADS-GAN | 0.052±0.01 | 0.054±0.02 | 0.039±0.02 | 0.036±0.01 | 0.041±0.01 |
| | CTGAN | 0.031±0.01 | **0.011**±**0.01** | 0.015±0.01 | **0.008**±**0.01** | 0.015±0.01 |
| | TVAE | 0.054±0.01 | 0.040±0.01 | 0.017±0.01 | 0.034±0.01 | 0.013±0.01 |
| | NFlows | 0.048±0.01 | 0.038±0.01 | 0.044±0.01 | 0.030±0.01 | **0.007**±**0.01** |
| Wasserstein Distance (Lower is Better) | SurvivalGAN | **0.153**±**0.02** | 0.228±0.11 | **0.441**±**0.12** | 0.420±0.29 | **5.861**±**3.01** |
| | PrivBayes | 1.080±0.27 | 0.101±0.01 | 0.828±0.02 | 0.146±0.01 | 15.718±4.39 |
| | ADS-GAN | 1.694±0.51 | 2.393±0.21 | 2.207±0.55 | 2.108±0.01 | 17.806±0.28 |
| | CTGAN | 0.851±0.07 | **0.032**±**0.01** | 0.621±0.05 | **0.019**±**0.01** | 9.343±3.81 |
| | TVAE | 1.967±0.05 | 1.430±0.20 | 0.892±0.30 | 1.902±0.10 | 9.421±2.85 |
| | NFlows | 1.701±0.19 | 0.325±0.05 | 2.086±0.03 | 0.169±0.04 | 14.649±0.76 |
| Optimism (Closer to Zero is Better) | SurvivalGAN | **-0.006**±**0.01** | **-0.012**±**0.03** | **-0.036**±**0.05** | -0.079±0.02 | -0.113±0.01 |
| | PrivBayes | **0.006**±**0.02** | -0.017±0.01 | 0.129±0.01 | **-0.001**±**0.01** | 0.384±0.01 |
| | ADS-GAN | 0.038±0.04 | 0.113±0.01 | 0.096±0.21 | 0.026±0.01 | 0.290±0.03 |
| | CTGAN | -0.061±0.03 | -0.334±0.02 | -0.099±0.02 | -0.386±0.04 | 0.183±0.08 |
| | TVAE | 0.064±0.01 | 0.079±0.01 | -0.115±0.03 | 0.024±0.01 | 0.140±0.03 |
| | NFlows | -0.183±0.05 | -0.236±0.06 | 0.091±0.02 | -0.398±0.10 | **0.003**±**0.01** |
| KM Divergence (Lower is Better) | SurvivalGAN | **0.011**±**0.01** | **0.031**±**0.02** | **0.079**±**0.01** | 0.080±0.02 | 0.121±0.01 |
| | PrivBayes | 0.035±0.01 | 0.034±0.01 | 0.134±0.01 | **0.007**±**0.01** | 0.383±0.01 |
| | ADS-GAN | 0.054±0.02 | 0.113±0.01 | 0.167±0.15 | 0.026±0.01 | 0.308±0.03 |
| | CTGAN | 0.073±0.02 | 0.334±0.02 | 0.099±0.02 | 0.386±0.04 | 0.133±0.04 |
| | TVAE | 0.064±0.01 | 0.083±0.01 | 0.115±0.03 | 0.024±0.01 | 0.168±0.03 |
| | NFlows | 0.185±0.05 | 0.238±0.06 | 0.100±0.01 | 0.392±0.10 | **0.054**±**0.01** |
| Short-sightedness (Closer to Zero is Better) | SurvivalGAN | 0.007±0.01 | 0.046±0.05 | 0.127±0.15 | 0.010±0.04 | **0.027**±**0.02** |
| | PrivBayes | 0.013±0.01 | 0.004±0.01 | 0.135±0.01 | **0.000**±**0.00** | ∗ |
| | ADS-GAN | 0.074±0.02 | 0.137±0.10 | 0.479±0.05 | 0.080±0.05 | ∗ |
| | CTGAN | **0.000**±**0.00** | **0.001**±**0.01** | 0.438±0.01 | **0.000**±**0.00** | 0.131±0.03 |
| | TVAE | **0.000**±**0.00** | 0.003±0.01 | 0.410±0.01 | **0.000**±**0.00** | 0.147±0.03 |
| | NFlows | 0.001±0.01 | 0.002±0.01 | **0.106**±**0.09** | 0.0001±0.01 | ∗ |

3. **Ablation study:** In Section 6.3, we perform an ablation study to demonstrate and quantify the importance of each component of SurvivalGAN. This provides insight into what each component of SurvivalGAN offers the overall model, and by using the new time and censoring metrics we are able to determine *how* SurvivalGAN fails when certain components are missing.

**Benchmarks.** We compare SurvivalGAN against the following benchmarks: generative adversarial networks for anonymization (**ADS-GAN**) (Yoon et al., 2020); conditional generative adversarial networks for tabular data (**CTGAN**) (Xu et al., 2019); variational autoencoder for tabular data (**TVAE**) (Xu et al., 2019); a variant of Bayesian Networks (**PrivBayes**) (Zhang et al., 2017); and Normalizing flows for tabular data (**NFlows**) (Papamakarios et al., 2021). For a fair comparison, we preprocess the data using our tabular encoder for all methods that are not specifically adapted to support tabular data (ADS-GAN, Normalizing flows).

**Datasets.** We test SurvivalGAN on a variety of medical survival analysis datasets. The datasets are: (1) ACTG 320 clinical trial dataset (**AIDS**) (Hammer et al., 1997); (2) Cambridge Urology Translational Research and Clinical Trials dataset for prostate cancer mortality in the UK (**CUTRACT**) (CUTRACT, 2019); (3) a private heart failure dataset (**PHEART**); (4) SEER dataset for prostate cancer mortality in the US (**SEER**) (SEER, 2019) and (5) The

Molecular Taxonomy of Breast Cancer International Consortium dataset (**METABRIC**) (Pereira et al., 2016). Details on each dataset can be found in Appendix C.1.

**Evaluation.** For each dataset, benchmark, and experimental setting, evaluations are performed using 5 different random seeds, and we report the mean and standard deviations of the desired metric. Further experimental details are provided in Appendix C and additional experiments in Appendix D. Code reproducing all experiments and implementing SurvivalGAN is publicly available at: `https://github.com/vanderschaarlab/survivalgan`.

### 6.1 Quality of Marginal Distributions

**Covariates.** To evaluate the quality of the covariates, in Table 1 we report the values of the Jensen-Shannon distance and Wasserstein distance. We do not expect SurvivalGAN to produce higher quality covariates than the baselines but we must confirm that they are not significantly worse. We see that SurvivalGAN provides robust coverage of the covariate space, as well as the benchmarks, often achieving the best or close to the best score.

**Time & Censoring.** To evaluate $p(t, E)$ we report the optimism, KM divergence, and short-sightedness metrics in Table 1. SurvivalGAN shows stable results across all datasets, consistently achieving the best value or close to

Table 2: Predictive performance of discriminative models trained with synthetic data. ∗: The discriminative models failed to train on the generated data. We provide the results when training with the original real data for reference.

| Metric | Method | AIDS | CUTRACT | PHEART | SEER | METABRIC |
|---|---|---|---|---|---|---|
| C-Index<br>(Higher is Better) | SurvivalGAN | **0.678±0.03** | **0.799±0.02** | **0.638±0.01** | **0.835±0.01** | **0.734±0.01** |
| | PrivBayes | 0.504±0.09 | 0.544±0.15 | 0.557±0.01 | 0.550±0.16 | 0.334±0.24 |
| | ADS-GAN | 0.541±0.06 | 0.607±0.10 | 0.565±0.03 | ∗ | 0.546±0.03 |
| | CTGAN | 0.546±0.08 | 0.791±0.02 | 0.612±0.01 | 0.830±0.01 | 0.499±0.04 |
| | TVAE | 0.564±0.04 | 0.779±0.03 | 0.621±0.02 | 0.808±0.02 | 0.711±0.02 |
| | NFlows | 0.522±0.08 | 0.730±0.07 | 0.566±0.02 | 0.776±0.04 | 0.482±0.01 |
| | Original | 0.742±0.02 | 0.826±0.01 | 0.668±0.01 | 0.856±0.01 | 0.706±0.02 |
| Brier Score<br>(Lower is Better) | SurvivalGAN | **0.057±0.01** | **0.084±0.01** | **0.181±0.01** | **0.023±0.01** | 0.179±0.01 |
| | PrivBayes | 0.058±0.01 | 0.100±0.02 | 0.209±0.02 | 0.025±0.02 | 0.562±0.31 |
| | ADS-GAN | 0.060±0.01 | 0.117±0.01 | 0.231±0.03 | ∗ | 0.260±0.05 |
| | CTGAN | 0.061±0.01 | 0.172±0.03 | 0.188±0.02 | 0.115±0.03 | 0.183±0.02 |
| | TVAE | 0.061±0.02 | 0.099±0.01 | 0.206±0.02 | 0.025±0.01 | **0.161±0.01** |
| | NFlows | 0.097±0.03 | 0.171±0.04 | 0.192±0.01 | 0.164±0.08 | 0.173±0.01 |
| | Original | 0.072±0.01 | 0.095±0.01 | 0.166±0.01 | 0.025±0.01 | 0.161±0.001 |

Table 3: Source-of-Gain Analysis on Multiple Datasets. ∗: The discriminative model failed to train on the generated data.

| Metric | Method | AIDS | CUTRACT | PHEART | SEER |
|---|---|---|---|---|---|
| C-Index<br>(Higher is Better) | SurvivalGAN | **0.723±0.02** | **0.804±0.01** | **0.644±0.01** | **0.834±0.01** |
| | w/o Time Regressor | 0.688±0.03 | 0.719±0.07 | 0.558±0.02 | 0.677±0.02 |
| | w/o Imbalanced Sampling | ∗ | 0.671±0.12 | 0.590±0.02 | 0.504±0.01 |
| | w/o Temporal Sampling | 0.713±0.04 | 0.792±0.01 | 0.614±0.01 | 0.636±0.13 |
| | w/o Cond. GAN | 0.714±0.04 | 0.655±0.13 | 0.563±0.04 | 0.573±0.12 |
| Brier Score<br>(Lower is Better) | SurvivalGAN | **0.066±0.01** | **0.083±0.01** | **0.176±0.01** | **0.022±0.01** |
| | w/o Time Regressor | 0.144±0.02 | 0.182±0.02 | 0.233±0.02 | 0.252±0.05 |
| | w/o Imbalanced Sampling | ∗ | 0.109±0.01 | 0.229±0.01 | 0.025±0.01 |
| | w/o Temporal Sampling | 0.152±0.01 | 0.113±0.01 | 0.223±0.01 | 0.153±0.11 |
| | w/o Cond. GAN | 0.187±0.01 | 0.247±0.06 | 0.234±0.04 | 0.224±0.11 |
| Optimism<br>(Closer to Zero is Better) | SurvivalGAN | **-0.006±0.01** | **-0.012±0.03** | **-0.036±0.05** | **-0.079±0.02** |
| | w/o Time Regressor | 0.067±0.01 | -0.025±0.15 | -0.070±0.06 | 0.123±0.01 |
| | w/o Imbalanced Sampling | 0.066±0.01 | 0.051±0.11 | 0.048±0.16 | 0.096±0.01 |
| | w/o Temporal Sampling | -0.284±0.07 | -0.266±0.01 | -0.116±0.07 | -0.644±0.03 |
| | w/o Cond. GAN | -0.309±0.03 | -0.369±0.01 | -0.130±0.04 | -0.445±0.11 |
| Short-sightedness<br>(Closer to Zero is Better) | SurvivalGAN | **0.007±0.01** | **0.046±0.05** | **0.127±0.15** | **0.010±0.04** |
| | w/o Time Regressor | 0.009±0.01 | 0.117±0.08 | 0.497±0.05 | 0.040±0.01 |
| | w/o Imbalanced Sampling | 0.008±0.01 | 0.064±0.06 | 0.228±0.14 | 0.132±0.12 |
| | w/o Temporal Sampling | 0.012±0.01 | 0.051±0.07 | 0.128±0.12 | 0.145±0.05 |
| | w/o Cond. GAN | 0.019±0.02 | 0.127±0.03 | 0.085±0.13 | 0.049±0.01 |

the best optimism and KM divergences. On the whole, the majority of baselines do not suffer significantly from short-sightedness. We see that SurvivalGAN is always close to the best value for a given dataset. Crucially, we see that SurvivalGAN is never extremely short-sighted. Whereas we see ADS-GAN, CTGAN and TVAE can suffer extreme failure cases measured by short-sightedness. This underlines one of our main novelties: good coverage of both covariate and temporal space by handling the censoring of the data, to robustly generate survival data.

## 6.2 Downstream Performance

To assess downstream performance, we train a set of discriminative models on the synthetic data and test them on the real data, known as the Train on Synthetic Test on Real (TSTR) metric (Esteban et al., 2017). The discriminative models fall into different categories of survival models: linear models (CoxPH) (Cox, 1972), gradient boosting (SurvivalXGBoost) (Barnwal et al., 2022), random forests (RandomSurvivalForest) (Ishwaran et al., 2008), and neural networks (DeepHit)

(Lee et al., 2018). For each experiment, we report the concordance index (**C-Index**) (Harrell et al., 1982), a standard metric for assessing the quality of the ranking in survival models, and the **Brier Score** (Brier et al., 1950), which measures the calibration of the probabilistic predictions. We evaluate the performance using 3-fold cross-validation on the real data for each discriminative model. The generative models are trained with training data, the predictive models are then trained on the synthetic data and tested on a held-out test set. Table 2 shows the downstream performance of the survival models. We select the best-performing predictive model for each metric and report its score. SurvivalGAN consistently leads to better-performing survival models, both in terms of prediction quality (C-Index) and calibration (Brier Score).

## 6.3 Sources of Gain

SurvivalGAN is characterized by: (1) the time-to-event/censoring regressor; (2) the imbalanced sampling with respect to time horizons and censoring and (3) the condi-
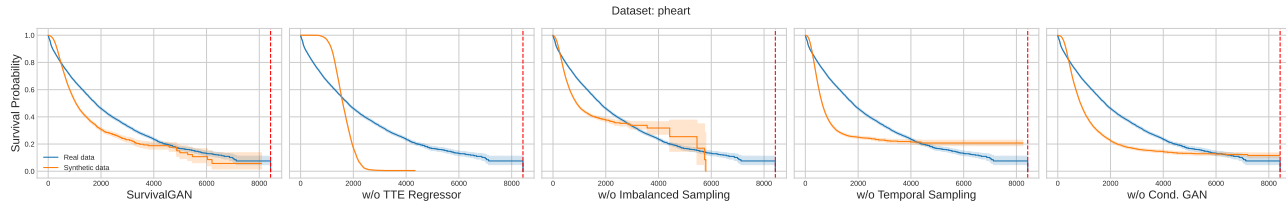
Figure 3: Sources of gain visualization for temporal quality using Kaplan-Meier plots. The sightedness is improved by the time regressor ($2^{nd}$ plot) and by the imbalanced sampling ($3^{rd}$ plot). The optimism is improved by the temporal sampling ($4^{th}$ plot) and by the conditional GAN ($5^{th}$ plot).

tional GAN used for the data generation. To examine the importance of each contribution, we apply the following modifications to SurvivalGAN: (1) with a standard time regressor (DATE) (Chapfuwa et al., 2018) instead of ours; (2) without the imbalanced sampling; (3) with imbalanced sampling, but only targeting the censoring instead of censoring and time horizon and (4) without the conditional GAN. We then evaluate downstream performance of predictive models. The ablation study was conducted in an in-distribution manner, leading to minor performance differences with Table 2. In addition to predictive performance, we also evaluate optimism and short-sightedness. We observe in Table 3 that all three components make significant contributions to improving the quality of the generated data.

**Insight.** The time regressor has an essential role in the quality of the rankings in the survival data (C-Index), as seen in the AIDS and CUTRACT datasets. It is also important with respect to short-sightedness: we see in particular for PHEART that, without a dedicated time regressor, the model is noticeably short-sighted. The imbalanced sampler plays a key role in the downstream models' ability to rank samples (measured by C-Index) for all datasets, but notably the large datasets, as seen in PHEART and SEER. Interestingly, without the temporal sampling or the conditional GAN, the model is severely pessimistic (negative optimism). Finally, without the conditional GAN, the model typically has a far worse Brier score across all datasets, showing that the conditional GAN significantly improves calibration.

We also visualize these sources of gain qualitatively using KM plots on the PHEART dataset in Figure 3. We see that all components of SurvivalGAN are required for optimal performance. This supports the idea that the imbalanced sampler and time regressor are crucial for the quality of the time values, while the conditional GAN is critical for temporal calibration.

## 7 CONCLUSION

We investigated the problem of generating synthetic survival data with censoring. We first formalized the problem, identifying three possible failure modes specific to the survival setting. We then introduced three new metrics based on

the survival functions to quantify these phenomena. On top of this, we introduced SurvivalGAN, a generative model that generates synthetic survival data. This is accomplished by incorporating censored and non-censored data into the training process, unlocking the use of abundant censored data. Additionally, the time-to-event/censoring data is generated in a more principled way, using a pre-trained survival function and time-to-event model, which permits future extensions to competing events. SurvivalGAN was tested on multiple medical datasets, generating more faithful data and leading to better downstream models than standard baseline generative methods.

**Limitations.** Currently SurvivalGAN is not able to address distribution shifts over time, where for example due to advances in medicine we might expect better survival rates in the future than we do now. In addition, it currently only operates in the static setting, with extensions to temporal data future work. Additionally, SurvivalGAN does not provide *guarantees* on privacy such as those in Jordon et al. (2018). We view these as exciting lines of future research.

## Acknowledgements

## References

Akrami, H., Joshi, A. A., Li, J., Aydöre, S., and Leahy, R. M. (2022). A robust variational autoencoder using beta divergence. *Knowledge-Based Systems*, 238:107886.

Alaa, A., van Breugel, B., Saveliev, E. S., and van der Schaar, M. (2022). How Faithful is your Synthetic Data? Sample-level Metrics for Evaluating and Auditing Generative Models. In *International Conference on Machine Learning*, pages 290–306. PMLR.

Ali, M. M., Malik, M. R., Ahmed, A. Y., Bashir, A. M., Mohamed, A., Abdi, A., and Obtel, M. (2022). Survival analysis of all critically ill patients with COVID-19 admitted to the main hospital in Mogadishu, Somalia, 30 March–12 June 2020: which interventions are proving effective in fragile states? *International Journal of Infectious Diseases*, 114:202–209.

Arsene, C. and Lisboa, P. (2007). Artificial Neural Networks Used in the Survival Analysis of Breast Cancer Patients: A Node-Negative Study. In *Outcome prediction in cancer*, pages 191–239. Elsevier.

Austin, P. C. (2012). Generating survival times to simulate cox proportional hazards models with time-varying covariates. *Statistics in Medicine*, 31(29):3946–3958.

Barnwal, A., Cho, H., and Hocking, T. (2022). Survival Regression with Accelerated Failure Time Model in XGBoost. *Journal of Computational and Graphical Statistics*, pages 1–25.

Bender, R., Augustin, T., and Blettner, M. (2005). Generating survival times to simulate Cox proportional hazards models. *Statistics in Medicine*, 24(11):1713–1723.

Benedetti, J. d., Oues, N., Wang, Z., Myles, P., and Tucker, A. (2020). Practical Lessons from Generating Synthetic Healthcare Data with Bayesian Networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 38–47. Springer.

Bretagnolle, J. and Huber, C. (1978). Estimation des densités: risque minimax. In *Séminaire de Probabilités XII*, pages 342–363. Springer.

Brier, G. W. et al. (1950). Verification of Forecasts Expressed in Terms of Probability. *Monthly weather review*, 78(1):1–3.

Chapfuwa, P., Tao, C., Li, C., Page, C., Goldstein, B., Duke, L. C., and Henao, R. (2018). Adversarial Time-to-Event Modeling. In *International Conference on Machine Learning*, pages 735–744. PMLR.

Chen, R. J., Lu, M. Y., Chen, T. Y., Williamson, D. F., and Mahmood, F. (2021). Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering*, 5(6):493–497.

Chen, T. and Guestrin, C. (2016). XGBoost. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM.

Cox, D. R. (1972). Regression Models and Life-Tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2):187–202.

Csiszár, I. and Körner, J. (2011). *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press.

CUTRACT (2019). UK, P. C. Prostate Cancer. https://prostatecanceruk.org/.

Danacica, D.-E. and Babucea, A.-G. (2010). Using Survival Analysis in Economics. *Survival*, 11:15.

De Capitani Di Vimercati, S., Foresti, S., Livraga, G., and Samarati, P. (2012). Data Privacy: Definitions and Techniques. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 20(06):793–817.

de Cos Juez, F. J., Nieto, P. G., Torres, J. M., and Castro, J. T. (2010). Analysis of lead times of metallic components in the aerospace industry through a supported vector machine model. *Mathematical and computer modelling*, 52(7-8):1177–1184.

Esteban, C., Hyland, S. L., and Rätsch, G. (2017). Real-valued (Medical) Time Series Generation with Recurrent Conditional GANs. *arXiv preprint arXiv:1706.02633*.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative Adversarial Nets. *Advances in Neural Information Processing Systems*, 27.

Gross, S. R., O'brien, B., Hu, C., and Kennedy, E. H. (2014). Rate of false conviction of criminal defendants who are sentenced to death. *Proceedings of the National Academy of Sciences*, 111(20):7230–7235.

Hammer, S. M., Squires, K. E., Hughes, M. D., Grimes, J. M., Demeter, L. M., Currier, J. S., Eron Jr, J. J., Feinberg, J. E., Balfour Jr, H. H., Deyton, L. R., et al. (1997). A Controlled Trial of Two Nucleoside Analogues Plus Indinavir in Persons with Human Immunodeficiency Virus Infection and CD4 Cell Counts of 200 per Cubic Millimeter or Less. *New England Journal of Medicine*, 337(11):725–733.

Harrell, F. E., Califf, R. M., Pryor, D. B., Lee, K. L., and Rosati, R. A. (1982). Evaluating the yield of medical tests. *JAMA*, 247(18):2543–2546.

Harsha, P., Jain, R., McAllester, D., and Radhakrishnan, J. (2007). The Communication Complexity of Correlation. In *IEEE Conference on Computational Complexity (CCC'07)*, volume 27, pages 10–23. IEEE.

Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. (2017). GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. *Advances in Neural Information Processing Systems*, 30.

Hinterstoisser, S., Pauly, O., Heibel, H., Martina, M., and Bokeloh, M. (2019). An Annotation Saved is an Annotation Earned: Using Fully Synthetic Training for Object

Instance Detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*.

Ho, J., Jain, A., and Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. *Advances in Neural Information Processing Systems*, 33:6840–6851.

Imrie, F., Cebere, B., McKinney, E. F., and van der Schaar, M. (2022). AutoPrognosis 2.0: Democratizing Diagnostic and Prognostic Modeling in Healthcare with Automated Machine Learning. *arXiv preprint arXiv:2210.12090*.

Ishwaran, H., Kogalur, U. B., Blackstone, E. H., and Lauer, M. S. (2008). Random survival forests. *The Annals of Applied Statistics*, 2(3):841–860.

Jain, P., Gyanchandani, M., and Khare, N. (2016). Big Data Privacy: A Technological Perspective and Review. *Journal of Big Data*, 3(1):1–25.

Jenkins, S. P. (2005). Survival analysis. *Unpublished manuscript, Institute for Social and Economic Research, University of Essex, Colchester, UK*, 42:54–56.

Jordon, J., Szpruch, L., Houssiau, F., Bottarelli, M., Cherubin, G., Maple, C., Cohen, S. N., and Weller, A. (2022). Synthetic data–what, why and how? *arXiv preprint arXiv:2205.03257*.

Jordon, J., Yoon, J., and van Der Schaar, M. (2018). PATE-GAN: Generating Synthetic Data with Differential Privacy Guarantees. In *International Conference on Learning Representations*.

Kaplan, E. L. and Meier, P. (1958). Nonparametric Estimation from Incomplete Observations. *Journal of the American statistical association*, 53(282):457–481.

Karras, T., Laine, S., and Aila, T. (2019). A Style-based Generator Architecture for Generative Adversarial Networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410.

Kaso, A. W., Agero, G., Hurissa, Z., Kaso, T., Ewune, H. A., Hareru, H. E., and Hailu, A. (2022). Survival analysis of COVID-19 patients in Ethiopia: A hospital-based study. *Plos one*, 17(5):e0268280.

Kent, S. L. (2010). Predicting Abolition: A Cross-National Survival Analysis of the Social and Political Determinants of Death Penalty Statutes. *International Criminal Justice Review*, 20(1):56–72.

Kingma, D. P. and Dhariwal, P. (2018). Glow: Generative Flow with Invertible 1x1 Convolutions. *Advances in Neural Information Processing Systems*, 31.

Kingma, D. P. and Welling, M. (2013). Auto-Encoding Variational Bayes. *arXiv preprint arXiv:1312.6114*.

LeClere, M. J. (2005). Preface Modeling Time to Event: Applications of Survival Analysis in Accounting, Economics and Finance. *Review of Accounting and Finance*.

Lee, C., Zame, W., Yoon, J., and van Der Schaar, M. (2018). DeepHit: A Deep Learning Approach to Survival Analysis with Competing Risks. In *Proceedings of the AAAI conference on artificial intelligence*.

Lee, E. T. and Go, O. T. (1997). Survival analysis in public health research. *Annual review of public health*, 18:105.

Liu, T., Qian, Z., Berrevoets, J., and van der Schaar, M. (2023). GOGGLE: Generative Modelling for Tabular Data by Learning Relational Structure. In *International Conference on Learning Representations*.

Lu, W., Yu, S., Liu, H., Suo, L., Tang, K., Hu, J., Shi, Y., and Hu, K. (2021). Survival Analysis and Risk Factors in COVID-19 Patients. *Disaster Medicine and Public Health Preparedness*, pages 1–6.

Luo, Y. and Lu, B.-L. (2018). EEG Data Augmentation for Emotion Recognition Using a Conditional Wasserstein GAN. In *2018 40th Annual International Conference of the IEEE engineering in Medicine and Biology Society*, pages 2535–2538. IEEE.

Machin, D., Cheung, Y. B., and Parmar, M. (2006). *Survival Analysis: A Practical Approach*. John Wiley & Sons.

Miscouridou, X., Perotte, A., Elhadad, N., and Ranganath, R. (2018). Deep Survival Analysis: Nonparametrics and Missingness. In *Machine Learning for Healthcare Conference*, pages 244–256. PMLR.

Naeem, M. F., Oh, S. J., Uh, Y., Choi, Y., and Yoo, J. (2020). Reliable Fidelity and Diversity Metrics for Generative Models. In *International Conference on Machine Learning*, pages 7176–7185. PMLR.

Niedermayer, D. (2008). An Introduction to Bayesian Networks and their Contemporary Applications. *Innovations in Bayesian networks: Theory and applications*, pages 117–130.

Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. (2016). WaveNet: A Generative Model for Raw Audio. *arXiv preprint arXiv:1609.03499*.

Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2021). Normalizing Flows for Probabilistic Modeling and Inference. *Journal of Machine Learning Research*, 22(57):1–64.

Pereira, B., Chin, S., and Rueda, O. e. a. (2016). The somatic mutation profiles of 2,433 breast cancers refine their genomic and transcriptomic landscapes. `https://www.nature.com/articles/ncomms11479`.

Perez, L. and Wang, J. (2017). The Effectiveness of Data Augmentation in Image Classification using Deep Learning. *arXiv preprint arXiv:1712.04621*.

Prakash, A., Boochoon, S., Brophy, M., Acuna, D., Cameracci, E., State, G., Shapira, O., and Birchfield, S. (2019). Structured Domain Randomization: Bridging the Reality

Gap by Context-Aware Synthetic Data. In *2019 International Conference on Robotics and Automation*, pages 7249–7255. IEEE.

Rajabi, A. and Garibay, O. O. (2022). TabFairGan: Fair Tabular Data Generation with Generative Adversarial Networks. *Machine Learning and Knowledge Extraction*, 4(2):488–501.

Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. (2022). Hierarchical Text-Conditional Image Generation with CLIP Latents. *arXiv preprint arXiv:2204.06125*.

Ranganath, R., Perotte, A., Elhadad, N., and Blei, D. (2016). Deep Survival Analysis. In *Machine Learning for Healthcare Conference*, pages 101–114. PMLR.

Reynolds, D. A. (2009). Gaussian Mixture Models. *Encyclopedia of biometrics*, 741(659-663).

Rezende, D. and Mohamed, S. (2015). Variational Inference with Normalizing Flows. In *International Conference on Machine Learning*, pages 1530–1538. PMLR.

Sajjadi, M. S., Bachem, O., Lucic, M., Bousquet, O., and Gelly, S. (2018). Assessing Generative Models via Precision and Recall. *Advances in Neural Information Processing Systems*, 31.

Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. (2016). Improved Techniques for Training GANs. *Advances in Neural Information Processing Systems*, 29.

Salinas-Escudero, G., Carrillo-Vega, M. F., Granados-García, V., Martínez-Valverde, S., Toledano-Toledano, F., and Garduño-Espinosa, J. (2020). A Survival Analysis of COVID-19 in the Mexican Population. *BMC public health*, 20(1):1–8.

SEER (2019). Surveillance, Epidemiology, and End Results (SEER) Program. `https://www.seer.cancer.gov/`.

Singh, R. and Mukhopadhyay, K. (2011). Survival Analysis in Clinical Trials: Basics and Must Know Areas. *Perspectives in clinical research*, 2(4):145.

Song, Y. and Ermon, S. (2019). Generative Modeling by Estimating Gradients of the Data Distribution. *Advances in Neural Information Processing Systems*, 32.

Sutherland, D. J., Tung, H.-Y., Strathmann, H., De, S., Ramdas, A., Smola, A., and Gretton, A. (2017). Generative Models and Model Criticism via Optimized Maximum Mean Discrepancy. In *International Conference on Learning Representations*.

Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., and Abbeel, P. (2017). Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 23–30. IEEE.

Tsybakov, A. B. (2009). Introduction to Nonparametric Estimation. *Springer Series in Statistics*.

Vahdat, A. and Kautz, J. (2020). NVAE: A Deep Hierarchical Variational Autoencoder. *Advances in Neural Information Processing Systems*, 33:19667–19679.

Vahdat, A., Kreis, K., and Kautz, J. (2021). Score-based Generative Modeling in Latent Space. *Advances in Neural Information Processing Systems*, 34:11287–11302.

van Breugel, B., Kyono, T., Berrevoets, J., and van der Schaar, M. (2021). DECAF: Generating fair Synthetic Data Using Causally-Aware Generative Networks. *Advances in Neural Information Processing Systems*, 34:22221–22233.

van der Maaten, L. and Hinton, G. (2008). Visualizing Data using t-SNE. *Journal of machine learning research*, 9(11).

Wang, Z., Myles, P., and Tucker, A. (2021). Generating and evaluating cross-sectional synthetic electronic healthcare data: Preserving data utility and patient privacy. *Computational Intelligence*, 37(2):819–851.

Xu, L., Skoularidou, M., Cuesta-Infante, A., and Veeramachaneni, K. (2019). Modeling Tabular Data using Conditional GAN. *Advances in Neural Information Processing Systems*, 32.

Xu, L. and Veeramachaneni, K. (2018). Synthesizing Tabular Data using Generative Adversarial Networks. *arXiv preprint arXiv:1811.11264*.

Yoon, J., Drumright, L. N., and van Der Schaar, M. (2020). Anonymization Through Data Synthesis Using Generative Adversarial Networks (ADS-GAN). *IEEE Journal of Biomedical and Health Informatics*, 24(8):2378–2388.

Zhang, J., Cormode, G., Procopiuc, C. M., Srivastava, D., and Xiao, X. (2017). PrivBayes: Private Data Release via Bayesian Networks. *ACM Trans. Database Syst.*, 42(4).

Zhou, X., Su, W., Liu, C., Jiao, Y., Zhao, X., and Huang, J. (2022). Deep Generative Survival Analysis: Nonparametric Estimation of Conditional Survival Function. *arXiv preprint arXiv:2205.09633*.

# Supplementary Material for:
# SurvivalGAN: Generating time-to-event Data for Survival Analysis

## A   BROADER IMPACT

**Applications.**   In general, there exist malicious applications of generative models. For example, when creating image or speech data, it is possible to make counterfeit images or recordings of an opponent, falsely damaging their reputation. Our paper is on survival data and we do not envision such malicious uses. Survival data is specific to the domain it is being applied to, which limits the possibility of using fake data for unethical purposes. As with all generative models, it is possible to reinforce biases in the training data. However, by providing a condition to the generator we are able to sample from underrepresented groups and our work makes progress in that respect.

Further, the use of SurvivalGAN allows us to train survival models without original data, removing real, sensitive human data from the supervised training process. This also makes it possible to generate more data quickly, helping to train survival models, which we have seen have positive societal impacts (Gross et al., 2014; Arsene and Lisboa, 2007; Danacica and Babucea, 2010).

**Datasets.**   In our experiments we use medical datasets. In general, these can contain sensitive information about participants. We note that we did not collect any data; all data was collected by external labs/medical researchers. All identifiable information has been removed from the datasets by the curators and permission was given by the subjects, making the datasets suitable for this paper.

## B   METRICS

In Section 4, we introduced three new metrics: optimism, short-sightedness, and Kaplan-Meier divergence. These metrics capture different nuances in generating data for survival analysis. In particular, we can see that we require three metrics because it is possible for one value to be zero, while the other two are non-zero. Figure 4 provides an illustration of this, showing the Kaplan-Meier plots for different situations and how the metrics differ. While KM divergence is enough to indicate when synthetic data is inadequate, it does not tell us *how*, for this we need the other two metrics because they target specific failure modes when modeling $p(t, E)$, giving us an interpretable meaning. Note that a KM divergence of zero means that optimism and short-sightedness must also both be zero.



Figure 4: Illustrative examples of Kaplan-Meier plots to show the need for all three metrics. Left: Optimism is zero, but short-sightedness and KM divergence are greater than zero. Middle: Short-sightedness is zero but optimism and KM divergence are greater than zero. Right: Optimism and short-sightedness are both zero, but KM divergence is greater than zero.

**Bounds on Optimism.**   Optimism is trivially bounded by -1 and 1 from its definition. Here we demonstrate that it is possible to bound the optimism by twice the total variation between the synthetic and real distributions $p_{\text{Syn}}(t)$ and $p_{\text{Real}}(t)$; which may be a tighter bound than -1 and 1 in certain situations. Recall that $p(t)$ is the probability density that the event of interest happens at time $t$, and the survival function is defined as $S(t) = \int_t^\infty p(t')dt'$. These can both be conditioned on covariate $\boldsymbol{x}$ which is omitted here for clarity. Differentiating this expression with respect to $t$ we obtain

$$\frac{dS}{dt} = -p(t).$$

Now we look at the definition of optimism

$$\text{Optimism} = \frac{1}{T} \int_0^T \big( S_{\text{Syn}}(t) - S_{\text{Real}}(t) \big) dt.$$

Applying integration by parts we get

$$\text{Optimism} = \left( S_{\text{Syn}}(T) - S_{\text{Real}}(T) \right) + \int_0^T \frac{t}{T} \big( p_{\text{Syn}}(t) - p_{\text{Real}}(t) \big) dt.$$

Rewriting the first bracketed term using the definition of the survival function we get

$$\text{Optimism} = \int_T^\infty 1 \times \big( p_{\text{Syn}}(t) - p_{\text{Real}}(t) \big) dt + \int_0^T \frac{t}{T} \times \big( p_{\text{Syn}}(t) - p_{\text{Real}}(t) \big) dt.$$

We then define $f(t)$ piecewise as

$$f(t) = \begin{cases} t/T, & \text{for } 0 \le t \le T \\ 1, & \text{for } t > T \end{cases}.$$

Giving the optimism as a single integral

$$\text{Optimism} = \int_0^\infty f(t) \big( p_{\text{Syn}}(t) - p_{\text{Real}}(t) \big) dt.$$

We then use the fact that $f(t) \le 1$ and $(p_{\text{Syn}}(t) - p_{\text{Real}}(t)) \le |p_{\text{Syn}}(t) - p_{\text{Real}}(t)|$ to conclude that

$$\text{Optimism} < 2 \int_0^\infty \frac{1}{2} |p_{\text{Syn}}(t) - p_{\text{Real}}(t)| dt.$$

The strict inequality comes from the fact that $f(t) < 1$ for $t < T$. The integral is the definition of the total-variation divergence. We obtain the lower bound by considering the negative of the above step, that is $(p_{\text{Syn}}(t) - p_{\text{Real}}(t)) \ge -|p_{\text{Syn}}(t) - p_{\text{Real}}(t)|$. Giving us

$$-2\mathcal{D}_{TV}(p_{\text{Real}}||p_{\text{Syn}}) < \text{Optimism} < 2\mathcal{D}_{TV}(p_{\text{Real}}||p_{\text{Syn}}). \tag{5}$$

This allows us to apply other known bounds between different probability distances and the total variation. For example, we can use Pinkser's inequality (Csiszár and Körner, 2011) to bound Optimism with the Kullback-Leibler divergence

$$-\sqrt{2\mathcal{D}_{KL}(p_{\text{Real}}||p_{\text{Syn}})} < \text{Optimism} < \sqrt{2\mathcal{D}_{KL}(p_{\text{Real}}||p_{\text{Syn}})}. \tag{6}$$

We can also apply the inequality of Bretagnolle and Huber (Bretagnolle and Huber, 1978; Tsybakov, 2009)

$$-2\sqrt{1 - \exp(-\mathcal{D}_{KL}(p_{\text{Real}}||p_{\text{Syn}}))} < \text{Optimism} < 2\sqrt{1 - \exp(-\mathcal{D}_{KL}(p_{\text{Real}}||p_{\text{Syn}}))}. \tag{7}$$

Another example is with the Hellinger distance (Harsha et al., 2007)

$$-2\sqrt{2}\mathcal{D}_H(p_{\text{Real}}||p_{\text{Syn}}) < \text{Optimism} < 2\sqrt{2}\mathcal{D}_H(p_{\text{Real}}||p_{\text{Syn}}). \tag{8}$$

The total-variation distance is symmetric. Therefore, despite the Kullback-Leibler divergence and Hellinger distance not being symmetric themselves, we are able to swap $p_{\text{Real}}$ and $p_{\text{Syn}}$ in Equations (6), (7) and (8) to obtain the tightest bounds. In specific situations there may be a closed-form solution for some of these divergences, allowing us to quickly establish bounds on the optimism.

# C  EXPERIMENTAL DETAILS

In this section we give full experimental details, giving the dataset descriptions and hyperparameters.

**Reproducibility.**  All hyperparameters are given in this section for reproducibility. Three of our datasets (SEER, AIDS and METABRIC) are public, making those experiments straightforward to run. It is possible to run all experiments with limited compute. Ours were run locally on a machine with 32GB RAM, Intel Core i7-6700 HQ, GeForce GTX 950M.

## C.1  Data Description

In Table 4 we provide details about the datasets used in our evaluation. Of the datasets, AIDS, SEER and METABRIC are public, CUTRACT and PHEART are licensed. AIDS contains people with HIV and SEER with Prostate Cancer. CUTRACT is owned by Cambridge Prostate Cancer (`https://cambridgeprostatecancer.com/`)[2], and focuses on patients with Prostate Cancer. PHEART consists of patients from 30 medical studies who have experienced heart failure. The Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) database is a Canada-UK Project which contains targeted sequencing data of primary breast cancer samples, this version contains samples 1,980 and 689 features (Pereira et al., 2016).

Table 4: Datasets used for evaluation.

| Dataset | No. instances | No. censored instances | No. features | Experiment label |
|---|---|---|---|---|
| ACTG 320 clinical trial dataset | 1151 | 1055 | 11 | AIDS |
| CUTRACT | 10086 | 8881 | 6 | CUTRACT |
| PHEART | 40409 | 25664 | 29 | PHEART |
| SEER prostate cancer | 171942 | 167568 | 6 | SEER |
| METABRIC | 1093 | 609 | 689 | METABRIC |

**Evaluation Horizons.**  The evaluation process is averaged over 5 time horizons, specific to each dataset. Given time-to-event $T$ in the training set, we select the 5 evaluation horizons, evenly spaced between $[T.min(), T.max()]$. Table 5 gives the evaluation horizons for each dataset.

Table 5: Time horizons used for evaluation by dataset (can represent days/months, depending on the datasets).

| Dataset | Horizon 1 | Horizon 2 | Horizon 3 | Horizon 4 | Horizon 5 |
|---|---|---|---|---|---|
| AIDS (Months) | 61.5 | 122. | 182.5 | 243. | 303.5 |
| CUTRACT (Days) | 1051. | 2082. | 3113. | 4144. | 5175. |
| PHEART (Days) | 1404.6 | 2809.3 | 4214. | 5618.6 | 7023.3 |
| SEER (Days) | 775. | 1550. | 2325. | 3100. | 3875. |
| METABRIC (Months) | 56.25 | 112.41 | 168.56 | 224.72 | 280.87 |

---

[2]Terms and Conditions:
`https://cambridgeprostatecancer.com/terms-privacy-policy-gdpr-cookies/`

## C.2   Hyperparameters

In Table 6, we present the full configuration of SurvivalGAN used in our experiments. Table 7 contains the hyperparameters used by the predictive models trained with synthetic data for downstream tasks. Finally, Table 8 details the hyperparameters used for the synthetic benchmarks.

Table 6: SurvivalGAN Hyperparameters by Component.

| Component | Parameter | Parameter Value |
|---|---|---|
| Survival Function | Model | Deephit |
| | No. Durations | 100 |
| | Batch Size | 100 |
| | No. Epochs | 2000 |
| | Learning Rate | $1 \times 10^{-3}$ |
| | Hidden Width | 300 |
| | $\alpha$ | 0.28 |
| | $\sigma$ | 0.38 |
| | Dropout Rate | 0.02 |
| | Patience | 20 |
| | Using Batch Normalization | True |
| Time-to-event Regressor | Model | XGBoostRegressor |
| | No. Estimators | 200 |
| | Depth | 5 |
| | Booster | Dart |
| | Tree Method | Histogram |
| Conditional GAN | Data Sampling Strategy | With Data Frequencies |
| | No. Iterations | 1500 |
| | Generator No. Hidden Layers | 3 |
| | Generator Hidden Width | 250 |
| | Generator Non-linearity | Tanh |
| | Generator Dropout Rate | 0.1 |
| | Discriminator No. Hidden Layers | 2 |
| | Discriminator Hidden Width | 250 |
| | Discriminator Non-linearity | Leaky ReLU |
| | Discriminator Dropout Rate | 0.1 |
| | Learning Rate | $1 \times 10^{-3}$ |
| | Weight Decay | $1 \times 10^{-3}$ |
| | Batch Size | 500 |
| | Gradient Penalty ($\lambda$) | 10 |
| | Encoder Max Clusters | 10 |

Table 7: Hyperparameters used for the baseline time-to-event benchmarks, used as downstream models.

| Method | Parameter | Parameter Value |
|---|---|---|
| CoxPH | Estimation Method | Breslow |
| | Penalizer | 0.0 |
| | $L^1$ Ratio | 0.0 |
| Weibull AFT | $\alpha$ | 0.05 |
| | Penalizer | 0.0 |
| | $L^1$ Ratio | 0.0 |
| SurvivalXGBoost | Objective | Survival: AFT |
| | Evaluation Metric | AFT Negative Log Likelihood |
| | AFT Loss Distribution | Normal |
| | AFT Loss Distribution Scale | 1.0 |
| | No. Estimators | 100 |
| | Column Subsample Ratio (by node) | 0.5 |
| | Maximum Depth | 8 |
| | Subsample Ratio | 0.5 |
| | Learning Rate | $5 \times 10^{-2}$ |
| | Minimum Child Weight | 50 |
| | Tree Method | Histogram |
| | Booster | Dart |
| RandomSurvivalForest | Max Depth | 3 |
| | No. Estimators | 100 |
| | Criterion | Gini |
| Deephit | No. Durations: | 1000 |
| | Batch Size | 100 |
| | Epochs | 2000 |
| | Learning Rate | $1 \times 10^{-3}$ |
| | Hidden Width | 300 |
| | $\alpha$ | 0.28 |
| | $\sigma$ | 0.38 |
| | Dropout Rate | 0.02 |
| | Patience | 20 |
| | Using Batch Normalization | True |
| DATE | Generator No. Hidden Layers | 2 |
| | Generator Hidden Width | 250 |
| | Generator Non-linearity | Leaky ReLU |
| | Generator No. Iterations | 1000 |
| | Generator Using Batch Normalization | False |
| | Generator Dropout Rate | 0.0 |
| | Generator Learning Rate | $2 \times 10^{-4}$ |
| | Generator Weight Decay | $1 \times 10^{-3}$ |
| | Generator Residual | True |
| | Discriminator No. Hidden Layers | 3 |
| | Discriminator Hidden Width | 300 |
| | Discriminator Non-linearity | Leaky ReLU |
| | Discriminator No. Iterations | 1 |
| | Discriminator Using Batch Normalization | False |
| | Discriminator Dropout Rate | 0.1 |
| | Discriminator Learning Rate | $2 \times 10^{-4}$ |
| | Discriminator Weight Decay | $1 \times 10^{-3}$ |
| | Patience | 10 |
| | Batch Size | 100 |

Table 8: Hyperparameters for the synthetic benchmarks.

| Method | Parameter | Parameter Value |
|---|---|---|
| PrivBayes | $\epsilon$ | 1.0 |
| | $\theta$ Usefulness | 4 |
| | $\epsilon$ Split | 0.3 |
| ADS-GAN | No. Iterations | 1500 |
| | Generator no. Hidden Layers | 3 |
| | Generator Hidden Width | 250 |
| | Generator Non-linearity | Tanh |
| | Generator Dropout Rate | 0.1 |
| | Discriminator No. Hidden Layers | 2 |
| | Discriminator Hidden Width | 250 |
| | Discriminator Non-linearity | Leaky ReLU |
| | Discriminator Dropout Rate | 0.1 |
| | Learning Rate | $1 \times 10^{-3}$ |
| | Weight Decay | $1 \times 10^{-3}$ |
| | Batch Size | 500 |
| | Gradient Penalty ($\lambda$) | 10 |
| | Identifiability Penalty | 0.1 |
| | Encoder Max Clusters | 10 |
| CTGAN | Embedding Width | 10 |
| | Generator No. Hidden Layers | 2 |
| | Generator Hidden Width | 256 |
| | Generator Learning Rate | $2 \times 10^{-4}$ |
| | Generator Decay | $1 \times 10^{-6}$ |
| | Discriminator No. Hidden Layers | 2 |
| | Discriminator Hidden Width | 256 |
| | Discriminator Learning Rate | $2 \times 10^{-4}$ |
| | Discriminator Decay | $1 \times 10^{-6}$ |
| | Batch Size | 500 |
| | Discriminator Steps | 1 |
| | No. Iterations | 300 |
| | Pac | 10 |
| TVAE | Embedding Width | 128 |
| | Encoder No. Hidden Layers | 2 |
| | Encoder Hidden Width | 128 |
| | Decoder No. Hidden Layers | 2 |
| | Decoder Hidden Width | 128 |
| | $L^2$ Scale | $1 \times 10^{-5}$ |
| | Batch Size | 500 |
| | No. Iterations | 300 |
| | Loss Factor | 2 |
| NFlows | No. Iterations | 500 |
| | No. Hidden Layers | 1 |
| | Hidden Width | 100 |
| | Batch Size | 100 |
| | No. Transform Blocks | 1 |
| | Dropout Rate | 0.1 |
| | No. Bins | 8 |
| | Tail Bound | 3 |
| | Learning Rate | $1 \times 10^{-3}$ |
| | Base Distribution | Standard Normal |
| | Linear Transform Type | Permutation |
| | Base Transform Type | Affine-Coupling |

# D  ADDITIONAL EXPERIMENTS

Here we present extended results from the experiments in the main paper. We first conduct further evaluation of the covariates. We then look at three further ablations of SurvivalGAN: using a dedicated censoring network; replacing the conditional GAN with a conditional VAE and replacing the GAN with a Gaussian mixture model. Finally, we provide qualitative results in the form of t-SNE and Kaplan-Meier plots in Appendix E.

## D.1  Statistical Metrics for the Covariates

In Table 1, we reported the Jensen-Shannon distance and Wasserstein distance of synthetic covariates compared to real covariates. We saw that SurvivalGAN generates data close to the real distribution, either similar to or better than the baselines. Here we extend this and record four further metrics:

1. The Precision (Sajjadi et al., 2018) - measures the rate by which the generative model synthesizes realistic-looking samples, higher is better.
2. The Recall (Sajjadi et al., 2018) - measures the fraction of real samples that are covered by the synthetic data, higher is better.
3. The Density (Naeem et al., 2020) - improves the precision metric and measures how many real-sample neighborhood spheres contain the generated data, higher is better.
4. The Coverage (Naeem et al., 2020) - improves the recall metric and reports the fraction of real samples whose neighborhoods contain at least one generated sample, higher is better.

We give these values in Table 9. We see that SurvivalGAN again performs on par with the baselines. This shows that SurvivalGAN does not generate worse covariates than the baselines, which is not the main aim of this work, but is crucial to generate a full synthetic survival dataset.

Table 9: Precision, recall, density, and coverage given the covariates manifolds, higher values are better. The SEER dataset evaluation failed due to memory limits. ∗: Not evaluated.

| Metric | Method | AIDS | CUTRACT | PHEART | METABRIC |
|---|---|---|---|---|---|
| Precision<br>(Higher is Better) | SurvivalGAN | 0.964±0.01 | 0.986±0.01 | 0.958±0.02 | 0.009±0.01 |
| | PrivBayes | 0.885±0.02 | 0.962±0.01 | 0.811±0.01 | ∗ |
| | ADS-GAN | 0.898±0.13 | **0.989±0.01** | 0.941±0.04 | 0.001±0.01 |
| | CTGAN | 0.878±0.04 | 0.976±0.01 | 0.929±0.03 | 0.006±0.01 |
| | TVAE | **0.970±0.01** | 0.960±0.02 | **0.987±0.01** | **0.014±0.01** |
| | NFlows | 0.880±0.03 | 0.864±0.02 | 0.582±0.02 | 0.004±0.01 |
| Recall<br>(Higher is Better) | SurvivalGAN | 0.911±0.02 | 0.714±0.10 | 0.689±0.06 | 0.698±0.04 |
| | PrivBayes | 0.968±0.01 | **0.981±0.01** | 0.975±0.01 | ∗ |
| | ADS-GAN | 0.829±0.10 | 0.621±0.10 | 0.489±0.08 | **0.982±0.02** |
| | CTGAN | 0.947±0.01 | 0.980±0.01 | 0.931±0.02 | 0.914±0.06 |
| | TVAE | 0.712±0.07 | 0.798±0.02 | 0.665±0.01 | 0.762±0.14 |
| | NFlows | **0.972±0.02** | 0.967±0.01 | **0.991±0.01** | 0.950±0.03 |
| Density<br>(Higher is Better) | SurvivalGAN | 1.018±0.06 | 0.976±0.04 | 0.902±0.05 | 0.009±0.01 |
| | PrivBayes | 0.719±0.03 | 0.876±0.01 | 0.557±0.02 | ∗ |
| | ADS-GAN | 0.991±0.25 | 1.001±0.10 | 0.889±0.10 | 0.001±0.01 |
| | CTGAN | 0.758±0.07 | 0.940±0.02 | 0.823±0.08 | 0.003±0.01 |
| | TVAE | **1.179±0.03** | **1.107±0.01** | **1.264±0.02** | **0.013±0.01** |
| | NFlows | 0.641±0.09 | 0.605±0.03 | 0.292±0.01 | 0.001±0.01 |
| Coverage<br>(Higher is Better) | SurvivalGAN | **0.919±0.02** | 0.505±0.06 | 0.546±0.01 | 0.034±0.02 |
| | PrivBayes | 0.834±0.04 | 0.903±0.01 | 0.685±0.01 | ∗ |
| | ADS-GAN | 0.792±0.09 | 0.472±0.02 | 0.493±0.05 | 0.002±0.01 |
| | CTGAN | 0.845±0.05 | **0.927±0.01** | **0.874±0.02** | 0.011±0.01 |
| | TVAE | 0.825±0.04 | 0.742±0.03 | 0.638±0.03 | **0.049±0.03** |
| | NFlows | 0.669±0.07 | 0.582±0.05 | 0.261±0.04 | 0.006±0.01 |

We additionally evaluate the quality of the covariates by looking at the negative log-likelihood of the synthetic covariates compared to that of the true covariates in Table 10. We see that SurvivalGAN typically matches the real data well, often the best or close to best, and never fails significantly, whereas the baselines occasionally contain noticeable failure cases, in particular CTGAN, TVAE and NFlows.

Table 10: Negative log-likelihood in the presence of the covariates. The closer to the real data the better. The values **closest** to the real data are given in bold. Extreme <u>failure</u> cases are underlined. ∗: Not evaluated.

| Source | AIDS $(/10^2)$ | CUTRACT $(/10^3)$ | PHEART $(/10^5)$ | SEER $(/10^4)$ | METABRIC $(/10^3)$ |
|---|---|---|---|---|---|
| Real data | 6.21 | 9.79 | 1.40 | 4.66 | 2.02 |
| SurvivalGAN | **6.37±0.72** | 11.24±4.51 | **1.40±0.25** | **4.41±2.21** | **1.35±0.14** |
| PrivBayes | 4.53±2.83 | **10.31±0.25** | **1.41±0.18** | **4.84±0.08** | ∗ |
| CTGAN | 10.04±1.04 | <u>47.72±1.38</u> | 2.03±0.02 | <u>99.95±16.05</u> | 1.19±0.19 |
| TVAE | <u>0.12±0.14</u> | <u>2.98±0.73</u> | 1.31±0.13 | <u>0.27±0.15</u> | 1.20±0.37 |
| NFlows | <u>26.42±3.78</u> | <u>55.33±1.82</u> | 1.97±0.39 | <u>93.79±30.93</u> | ∗ |

## D.2 Censoring network

As a further ablation, we aim to investigate if a dedicated censoring network would improve the quality of SurvivalGAN. We design the following experiment: We use an XGBoost classifier - denoted Censoring network - to predict the censored/not censored status, based on the covariates. We keep our mechanisms in place: unbalanced time/censoring sampling, and time-to-event/censoring regression. The results are given in Tables 11, 12 and 13. We see that on the whole, SurvivalGAN performs better without the censoring network.

Table 11: Censoring network training performance. We evaluate the classifier using only the real data and we report the AUROC. We observe that the classifier has a good performance for distinguishing the classes, on the evaluation datasets.

| Dataset | AUROC |
|---|---|
| AIDS | 0.719±0.037 |
| CUTRACT | 0.753±0.004 |
| PHEART | 0.719±0.004 |
| SEER | 0.837±0.003 |
| METABRIC | 0.717±0.015 |

Table 12: Number of censored/not censored rows, from the real data, SurvivalGAN, and the censoring network. We want the numbers from the generative models to be as close as possible to those from the real data. While the predictive performance is good overall, there are scenarios like for the AIDS or CUTRACT datasets, where the number of non-censored synthetic subjects from the Censoring Network is too low.

| Dataset | Real data | | SurvivalGAN | | Censoring network | |
|---|---|---|---|---|---|---|
| | Censored | Event | Censored | Event | Censored | Event |
| AIDS | 1055 | 96 | 1045 | 106 | 1134 | 17 |
| CUTRACT | 8881 | 1205 | 8868 | 1218 | 9926 | 160 |
| MAGGIC | 25664 | 14745 | 26019 | 14390 | 27282 | 13127 |
| SEER | 167568 | 4374 | 168866 | 3076 | 132260 | 4872 |
| METABRIC | 609 | 484 | 817 | 276 | 738 | 355 |

Table 13: Predictive performance for SurvivalGAN with and without the Censoring network. C-Index is better if higher and Brier score is better if lower. We see SurvivalGAN performs better without the censoring network.

| Dataset | SurvivalGAN | | Censoring network | |
|---|---|---|---|---|
| | C-Index | Brier score | C-Index | Brier score |
| AIDS | 0.723±0.020 | **0.066±0.010** | **0.730±0.010** | 0.067±0.001 |
| CUTRACT | **0.804±0.010** | **0.083±0.010** | 0.770±0.010 | 0.100±0.001 |
| PHEART | **0.644±0.010** | **0.176±0.010** | 0.640±0.005 | 0.179±0.004 |
| SEER | **0.834±0.010** | **0.022±0.010** | 0.774±0.001 | 0.032±0.002 |
| METABRIC | **0.719±0.020** | **0.200±0.002** | 0.710±0.005 | 0.210±0.005 |

## D.3 SurvivalGAN vs. SurvivalVAE

Our method can be adapted to other architectures as well to generate covariates. In this section, we perform another ablation by analyzing the performance of the synthetic data when using a variational autoencoder (**SurvivalVAE**) instead of a GAN. For the experiment, we keep the same additional mechanisms in-place: imbalanced sampling around time and censoring, and the time-to-event/censoring regression. Table 14 contains the predictive performance of models trained on the synthetic data. We see that SurvivalGAN outperforms SurvivalVAE on the majority of datasets.

Table 14: Predictive performance for models trained with synthetic data from SurvivalGAN vs. SurvivalVAE. C-Index is better if higher and Brier score is better if lower. SurvivalGAN tends to generate better quality data.

| Dataset | Method | C-Index | Brier score |
|---|---|---|---|
| AIDS | SurvivalVAE | 0.638±0.020 | 0.058±0.000 |
| | SurvivalGAN | **0.678±0.030** | **0.057±0.010** |
| CUTRACT | SurvivalVAE | 0.791±0.010 | 0.103±0.003 |
| | SurvivalGAN | **0.799±0.020** | **0.084±0.010** |
| PHEART | SurvivalVAE | 0.600±0.001 | 0.206±0.002 |
| | SurvivalGAN | **0.638±0.010** | **0.181±0.010** |
| SEER | SurvivalVAE | 0.609±0.010 | 0.024±0.010 |
| | SurvivalGAN | **0.835±0.010** | **0.023±0.010** |
| METABRIC | SurvivalVAE | 0.724±0.010 | 0.191±0.002 |
| | SurvivalGAN | **0.734±0.010** | **0.189±0.010** |

## D.4 SurvivalGAN vs. SurvivalGMM

We test SurvivalGAN against the simplest possible generative model for the covariates, a Gaussian Mixture Model (SurvivalGMM). This uses the same Gaussian Mixture Model that is used in the tabular encoder and class encoder with 100 mixture components (see Section 5.1 for information on these), but now to generate covariates instead of using the GAN. We provide the downstream performances in Table 15. We see that using a GAN performs better than using the GMM.

Table 15: Predictive performance for models trained with synthetic data from SurvivalGAN vs. SurvivalGMM. C-Index is better if higher and Brier score is better if lower. We see SurvivalGAN is significantly better than SurvivalGMM.

| Dataset | Method | C-Index | Brier score |
|---|---|---|---|
| AIDS | SurvivalGMM | 0.510±0.157 | 0.061±0.002 |
| | SurvivalGAN | **0.678±0.030** | **0.057±0.010** |
| CUTRACT | SurvivalGMM | 0.780±0.008 | 0.089±0.001 |
| | SurvivalGAN | **0.799±0.020** | **0.084±0.010** |
| PHEART | SurvivalGMM | 0.627±0.002 | 0.205±0.001 |
| | SurvivalGAN | **0.638±0.010** | **0.181±0.010** |
| SEER | SurvivalGMM | 0.662±0.017 | 0.024±0.010 |
| | SurvivalGAN | **0.835±0.010** | **0.023±0.010** |
| METABRIC | SurvivalGMM | 0.564±0.013 | 0.282±0.036 |
| | SurvivalGAN | **0.734±0.010** | **0.189±0.010** |

# E QUALITATIVE RESULTS

Here we provide qualitative results in the form of t-SNE plots for the covariates (van der Maaten and Hinton, 2008) and Kaplan-Meier plots for the time and event (Kaplan and Meier, 1958).

## E.1 Downstream Predictive Models

Figure 5 presents the Kaplan-Meier plots for the time-to-event models (the downstream models). The observed trend is visible in the datasets, leading to over-optimistic or over-pessimistic time-to-event values, and it supports the need for a reliable method to overcome the time-to-event/censoring problem.



Figure 5: Kaplan-Meier plots of reference time-to-event models. The first five survival-function-based methods tend to be over-optimistic, while the last model is over-pessimistic.

## E.2   Data Fidelity and Diversity

Figure 6, includes the t-SNE plots for covariate coverage and the Kaplan-Meier visualizations for temporal fidelity, using all datasets. Qualitatively, we see SurvivalGAN is robust in generating the covariate and temporal distributions across all datasets. The t-SNE plots show the covariates typically cover the data distribution at least as well as the other baselines. More importantly, the KM plots show the ability to model $p(t, E)$ is significantly better for SurvivalGAN than the baselines.



Figure 6: Data diversity visualization for all datasets, from a single random seed. For each dataset, the 1st row contains the t-SNE plots on the covariates, and the 2nd row contains the Kaplan-Meier plots for time and censoring. Each column provides the visualization for each of the available benchmarks.

## E.3  Sources of Gain

Figure 3 reports the sources of gain visualizations for all datasets (apart from METABRIC) with t-SNE for the covariates and KM plots for the time/censoring. We observe that the conditional GAN has an important impact on the quality of the covariates (final column of the t-SNE plots), while all the components contribute to the temporal calibration and sightedness (KM plots).
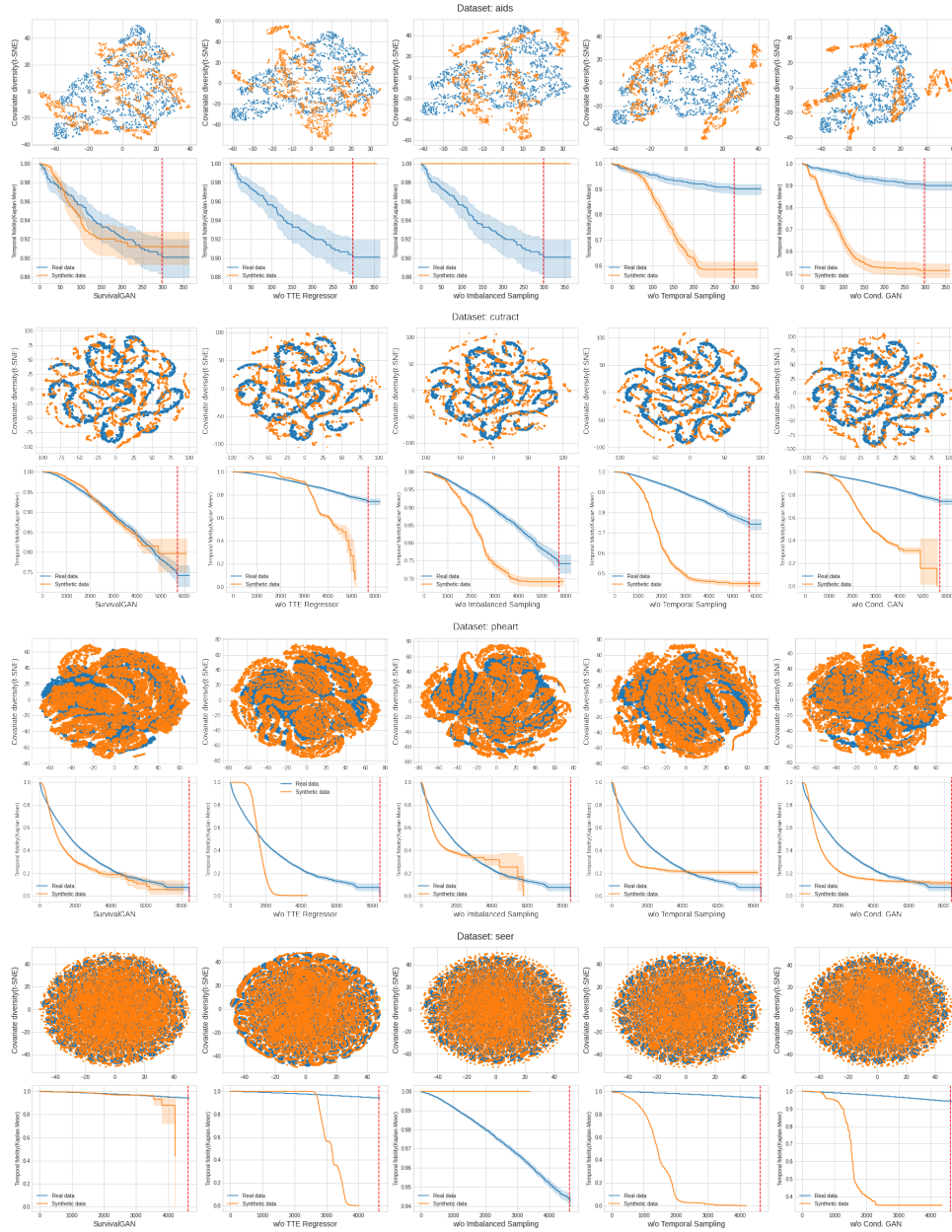


Figure 7: Sources of gain visualizations using t-SNE and Kaplan-Meier plots. For each dataset, the 1st row contains the t-SNE plots of the covariates, and the 2nd row contains the Kaplan-Meier plots for time and censoring. Each column corresponds to a given ablation scenario (certain component missing).

## E.4    SurvivalGAN vs. SurvivalVAE

Finally, Figure 8 shows the qualitative differences between SurvivalGAN and SurvivalVAE. The biggest difference is in the Kaplan-Meier plots, showing SurvivalGAN models $p(t, E)$ more faithfully than SurvivalVAE. We include CTGAN and TVAE as well to show that generally VAE based models are more over-optimistic than GAN based models.
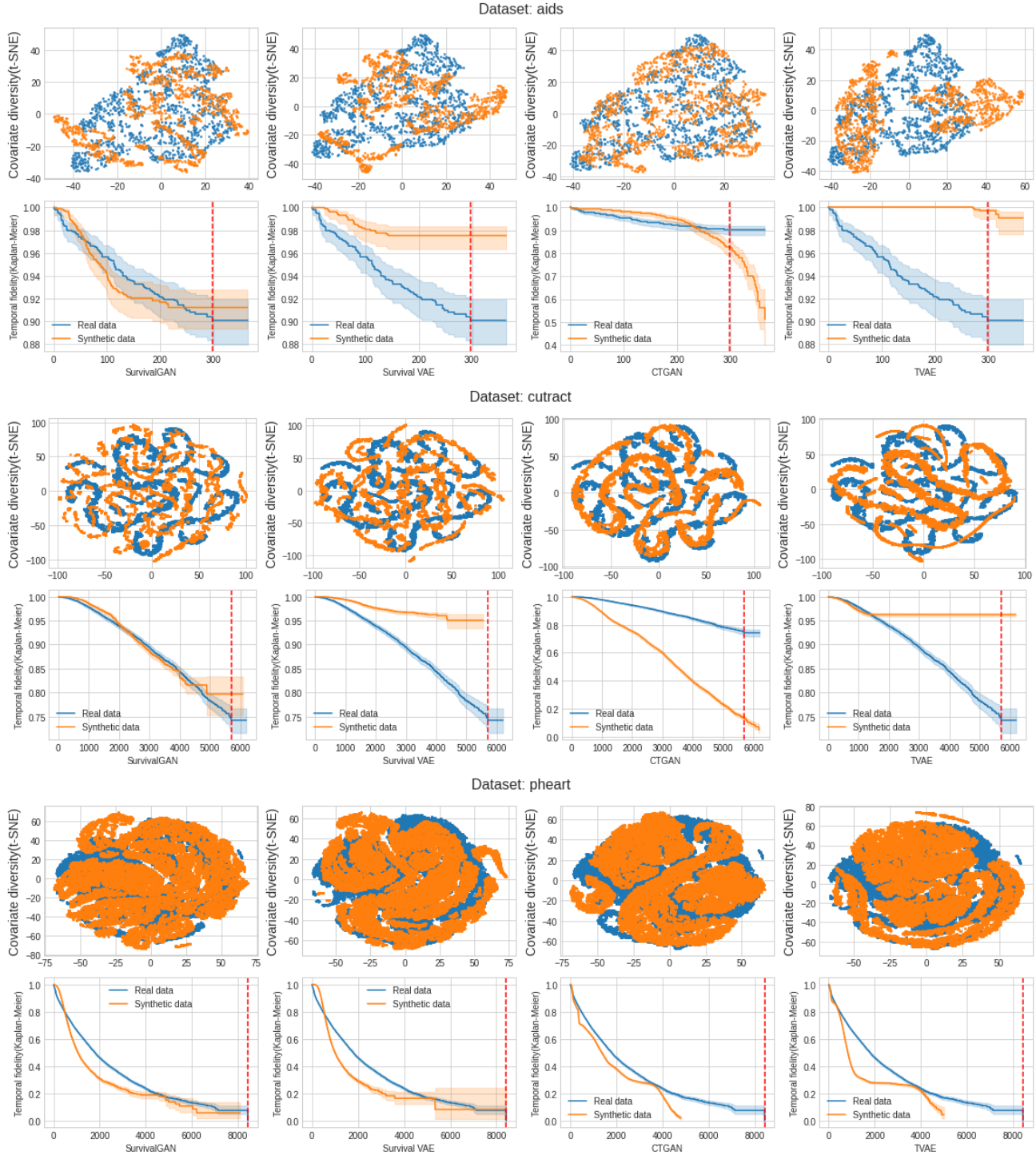


Figure 8: Data diversity visualization for SurvivalVAE, from a single random seed. For each dataset, the 1st row contains the t-SNE plots of the covariates, and the 2nd row contains the Kaplan-Meier plots for time and censoring. Each column provides the visualization for each of the 4 benchmarks.