

SurvLIME: A method for explaining machine learning survival models

Maxim S. Kovalev, Lev V. Utkin and Ernest M. Kasimov
Peter the Great St.Petersburg Polytechnic University (SPbPU)
St.Petersburg, Russia

e-mail: lev.utkin@gmail.com, maxkovalev03@gmail.com, kasimov.ernest@gmail.com

Abstract

A new method called SurvLIME for explaining machine learning survival models is proposed. It can be viewed as an extension or modification of the well-known method LIME. The main idea behind the proposed method is to apply the Cox proportional hazards model to approximate the survival model at the local area around a test example. The Cox model is used because it considers a linear combination of the example covariates such that coefficients of the covariates can be regarded as quantitative impacts on the prediction. Another idea is to approximate cumulative hazard functions of the explained model and the Cox model by using a set of perturbed points in a local area around the point of interest. The method is reduced to solving an unconstrained convex optimization problem. A lot of numerical experiments demonstrate the SurvLIME efficiency.

Keywords: interpretable model, explainable AI, survival analysis, censored data, convex optimization, the Cox model.

1 Introduction

Many complex problems in various applications are solved by means of deep machine learning models, in particular deep neural networks, at the present time. One of the demonstrative examples is the disease diagnosis by the models on the basis of medical images or another medical information. At the same time, deep learning models often work as black-box models such that details of their functioning are often completely unknown. It is difficult to explain in this case how a certain result or decision is achieved. As a result, the machine learning models meet some difficulties in their incorporating into many important applications, for example, into medicine, where doctors have to have an explanation of a stated diagnosis in order to choose a corresponding treatment. The lack of the explanation elements in many machine learning models has motivated development of many methods which could interpret or explain the deep machine learning algorithm predictions and understand the decision-making process or the key factors involved in the decision [4, 18, 35, 36].

The methods explaining the black-box machine learning models can be divided into two main groups: local methods which derive explanation locally around a test example; global methods which try to explain the overall behavior of the model. A key component of explanations for models is the contribution of individual input features. It is assumed that a prediction is explained when every feature is assigned by some number quantified its impact on the prediction. One of the first local explanation methods is the Local Interpretable Model-agnostic Explanations (LIME) [43], which uses simple and easily understandable linear models to locally approximate the predictions of black-box

models. The main intuition of the LIME is that the explanation may be derived locally from a set of synthetic examples generated randomly in the neighborhood of the example to be explained such that every synthetic example has a weight according to its proximity to the explained example. Moreover, the method uses simple and easily understandable models like decision rules or linear models to locally or globally approximate the predictions of black-box models. The method is agnostic to the black-box model. This means that any details of the black-box model are unknown. Only its input and the corresponding output are used for training the explanation model. It is important to mention a work [16], which provides a thorough theoretical analysis of the LIME.

The LIME as well as other methods have been successfully applied to many machine learning models for explanation. However, to the best of our knowledge, there is a large class of models for which there are no explanation methods. These models are applied to problems which take into account survival aspects of applications. Survival analysis as a basis for these models can be regarded as a fundamental tool which is used in many applied areas especially in medicine. The survival models can be divided into three parts: parametric, nonparametric and semiparametric [23, 56]. Most machine learning models are based on nonparametric and semiparametric survival models. One of the popular regression model for the analysis of survival data is the well-known Cox proportional hazards model, which is a semi-parametric model that calculates effects of observed covariates on the risk of an event occurring, for example, the death or failure [11]. The proportional hazards assumption in the Cox model means that different examples (patients) have hazard functions that are proportional, i.e., the ratio of the hazard functions for two examples with different prognostic factors or covariates is a constant and does not vary with time. The model assumes that a patient’s log-risk of failure is a linear combination of the example covariates. This is a very important peculiarity of the Cox model, which will be used below in explanation models.

A lot of machine learning implementations of survival analysis models have been developed [56] such that most implementations (random survival forests, deep neural networks) can be regarded as black-box models. Therefore, the problem of the survival analysis result explanation is topical. However, in contrast to other machine learning models, one of the main difficulties of the survival model explanation is that the result of most survival models is a time-dependent function (the survival function, the hazard function, the cumulative hazard function, etc.). This implies that many available explanation methods like LIME cannot be applied to survival models. In order to cope with this difficulty, a new method called SurvLIME (Survival LIME) for explaining machine learning survival models is proposed, which can be viewed as an extension or modification of LIME. The main idea of the proposed method is to apply the Cox proportional hazards model to approximate the survival model at a local area around a test example. It has been mentioned that the Cox model considers a linear combination of the example covariates. Moreover, it is important that the covariates as well as their combination do not depend on time. Therefore, coefficients of the covariates can be regarded as quantitative impacts on the prediction. However, we approximate not a point-valued black-box model prediction, but functions, for example, the cumulative hazard function (CHF). In accordance with the proposed explanation method, synthetic examples are randomly generated around the explainable example, and the CHF is calculated for every synthetic example by means of the black-box survival model. Simultaneously, we write the CHF corresponding to the approximating Cox model as a function of the coefficients of interest. By writing the distance between the CHF provided by the black-box survival model and the CHF of the approximating Cox model, we construct an unconstrained convex optimization problem for computing the coefficients of covariates. Numerical results using synthetic and real data illustrate the proposed method.

The paper is organized as follows. Related work can be found in Section 2. A short description of basic concepts of survival analysis, including the Cox model, is given in Section 3. Basic ideas of the

method LIME are briefly considered in Section 4. Section 5 provides a description of the proposed SurvLIME and its basic ideas. A formal derivation of the convex optimization problem for determining important features and a scheme of an algorithm implementing SurvLIME can be found in Section 6. Numerical experiments with synthetic data are provided in Section 7. Similar numerical experiments with real data are given in Section 8. Concluding remarks are provided in Section 9.

2 Related work

Local explanation methods. A lot of methods have been developed to locally explain black-box models. Along with the original LIME [43], many its modifications have been proposed due to success and simplicity of the method, for example, ALIME [47], NormLIME [2], DLIME [62], Anchor LIME [44], LIME-SUP [24], LIME-Aleph [40], GraphLIME [25]. Another very popular method is the SHAP [49] which takes a game-theoretic approach for optimizing a regression loss function based on Shapley values [33]. Alternative methods are influence functions [31], a multiple hypothesis testing framework [9], and many other methods.

An increasingly important family of methods are based on counterfactual explanations [54], which try to explain what to do in order to achieve a desired outcome by means of finding changes to some features of an explainable input example such that the resulting data point called counterfactual has a different specified prediction than the original input. Due to intuitive and human-friendly explanations provided by this family of methods, it is extended very quickly [17, 22, 32, 52]. Counterfactual modifications of LIME have been also proposed by Ramon et al. [41] and White and Garcez [57].

Many aforementioned explanation methods starting from LIME [43] are based on perturbation techniques [14, 15, 39, 53]. These methods assume that contribution of a feature can be determined by measuring how prediction score changes when the feature is altered [12]. One of the advantages of perturbation techniques is that they can be applied to a black-box model without any need to access the internal structure of the model. A possible disadvantage of perturbation technique is its computational complexity when perturbed input examples are of the high dimensionality.

Descriptions of many explanation methods and various approaches, their critical reviews can be found in survey papers [1, 3, 10, 18, 45].

It should be noted that most explanation methods deal with the point-valued predictions produced by black-box models. We mean under point-valued predictions some finite set of possible model outcomes, for example, classes of examples. A main problem of the considered survival models is that their outcome is a function. Therefore, we try to propose a new approach, which uses LIME as a possible tool for its implementing, dealing with CHF's as the model outcomes.

Machine learning models in survival analysis. A review of survival analysis methods is presented by Wang et al. [56]. The Cox model is a very powerful and popular method for dealing with survival data. Therefore, a lot of approaches modifying the Cox model have been proposed last decades. In particular, Tibshirani [51] proposed a modification of the Cox model based on the Lasso method in order to take into account a high dimensionality of survival data. Following this paper, several modifications of the Lasso methods for the Cox model were introduced [27, 30, 50, 59, 63]. In order to relax the linear relationship assumption accepted in the Cox model, a simple neural network as well as the deep neural networks were proposed by several authors [13, 19, 28, 42, 64]. The SVM approach to survival analysis has been also studied by several authors [6, 29, 48, 58]. It turned out that the random survival forests (RSFs) became a very powerful, efficient and popular tool for the survival analysis. Therefore, this tool and its modifications, which can be regarded as extensions of the standard random forest [8] on survival data, were proposed and investigated in many papers, for

example, in [26, 34, 37, 38, 46, 55, 60, 61], in order to take into account the limited survival data.

Most of the above models dealing with survival data can be regarded as black-box models and should be explained. However, only the Cox model has a simple explanation due to its linear relationship between covariates. Therefore, it can be used to approximate more powerful models, including survival deep neural networks and RSFs, in order to explain predictions of these models.

3 Some elements of survival analysis

3.1 Basic concepts

In survival analysis, an example (patient) i is represented by a triplet $(\mathbf{x}_i, \delta_i, T_i)$, where $\mathbf{x}_i^T = (x_{i1}, \dots, x_{id})$ is the vector of the patient parameters (characteristics) or the vector of the example features; T_i is time to event of the example. If the event of interest is observed, then T_i corresponds to the time between baseline time and the time of event happening, in this case $\delta_i = 1$, and we have an uncensored observation. If the example event is not observed and its time to event is greater than the observation time, then T_i corresponds to the time between baseline time and end of the observation, and the event indicator is $\delta_i = 0$, and we have a censored observation. Suppose a training set D consists of n triplets $(\mathbf{x}_i, \delta_i, T_i)$, $i = 1, \dots, n$. The goal of survival analysis is to estimate the time to the event of interest T for a new example (patient) with feature vector denoted by \mathbf{x} by using the training set D .

The survival and hazard functions are key concepts in survival analysis for describing the distribution of event times. The survival function denoted by $S(t)$ as a function of time t is the probability of surviving up to that time, i.e., $S(t) = \Pr\{T > t\}$. The hazard function $h(t)$ is the rate of event at time t given that no event occurred before time t , i.e.,

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{\Pr\{t \leq T \leq t + \Delta t | T \geq t\}}{\Delta t} = \frac{f(t)}{S(t)}, \quad (1)$$

where $f(t)$ is the density function of the event of interest.

By using the fact that the density function can be expressed through the survival function as

$$f(t) = -\frac{dS(t)}{dt}, \quad (2)$$

we can write the following expression for the hazard rate:

$$h(t) = -\frac{d}{dt} \ln S(t). \quad (3)$$

Another important concept in survival analysis is the CHF $H(t)$, which is defined as the integral of the hazard function $h(t)$ and can be interpreted as the probability of an event at time t given survival until time t , i.e.,

$$H(t) = \int_{-\infty}^t h(x) dx. \quad (4)$$

The survival function is determined through the hazard function and through the CHF as follows:

$$S(t) = \exp\left(-\int_0^t h(z) dz\right) = \exp(-H(t)). \quad (5)$$

The dependence of the above functions on a feature vector \mathbf{x} is omitted for short.

To compare survival models, the C-index proposed by Harrell et al. [21] is used. It estimates how good a survival model is at ranking survival times. It estimates the probability that, in a randomly selected pair of examples, the example that fails first had a worst predicted outcome. In fact, this is the probability that the event times of a pair of examples are correctly ranking.

3.2 The Cox model

Let us consider main concepts of the Cox proportional hazards model, [23]. According to the model, the hazard function at time t given predictor values \mathbf{x} is defined as

$$h(t|\mathbf{x}, \mathbf{b}) = h_0(t)\Psi(\mathbf{x}, \mathbf{b}) = h_0(t) \exp(\psi(\mathbf{x}, \mathbf{b})). \quad (6)$$

Here $h_0(t)$ is a baseline hazard function which does not depend on the vector \mathbf{x} and the vector \mathbf{b} ; $\Psi(\mathbf{x}, \mathbf{b})$ is the covariate effect or the risk function; $\mathbf{b}^T = (b_1, \dots, b_m)$ is an unknown vector of regression coefficients or parameters. It can be seen from the above expression for the hazard function that the reparametrization $\Psi(\mathbf{x}, \mathbf{b}) = \exp(\psi(\mathbf{x}, \mathbf{b}))$ is used in the Cox model. The function $\psi(\mathbf{x}, \mathbf{b})$ in the model is linear, i.e.,

$$\psi(\mathbf{x}, \mathbf{b}) = \mathbf{b}^T \mathbf{x} = \sum_{k=1}^m b_k x_k. \quad (7)$$

In the framework of the Cox model, the survival function $S(t|\mathbf{x}, \mathbf{b})$ is computed as

$$S(t|\mathbf{x}, \mathbf{b}) = \exp(-H_0(t) \exp(\psi(\mathbf{x}, \mathbf{b}))) = (S_0(t))^{\exp(\psi(\mathbf{x}, \mathbf{b}))}. \quad (8)$$

Here $H_0(t)$ is the cumulative baseline hazard function; $S_0(t)$ is the baseline survival function. It is important to note that functions $H_0(t)$ and $S_0(t)$ do not depend on \mathbf{x} and \mathbf{b} .

The partial likelihood in this case is defined as follows:

$$L(\mathbf{b}) = \prod_{j=1}^n \left[\frac{\exp(\psi(\mathbf{x}_j, \mathbf{b}))}{\sum_{i \in R_j} \exp(\psi(\mathbf{x}_i, \mathbf{b}))} \right]^{\delta_j}. \quad (9)$$

Here R_j is the set of patients who are at risk at time t_j . The term ‘‘at risk at time t ’’ means patients who die at time t or later.

4 LIME

Before studying the LIME modification for survival data, this method is briefly considered below.

LIME proposes to approximate a black-box model denoted as f with a simple function g in the vicinity of the point of interest \mathbf{x} , whose prediction by means of f has to be explained, under condition that the approximation function g belongs to a set of explanation models G , for example, linear models. In order to construct the function g in accordance with LIME, a new dataset consisting of perturbed samples is generated, and predictions corresponding to the perturbed samples are obtained by means of the explained model. New samples are assigned by weights $w_{\mathbf{x}}$ in accordance with their proximity to the point of interest \mathbf{x} by using a distance metric, for example, the Euclidean distance.

An explanation (local surrogate) model is trained on new generated samples by solving the following optimization problem:

$$\arg \min_{g \in G} L(f, g, w_{\mathbf{x}}) + \Phi(g). \quad (10)$$

Here L is a loss function, for example, mean squared error, which measures how the explanation is close to the prediction of the black-box model; $\Phi(g)$ is the model complexity.

A local linear model is the result of the original LIME. As a result, the prediction is explained by analyzing coefficients of the local linear model.

5 A basic sketch of SurvLIME

Suppose that there are available a training set D and an black-box model. For every new example \mathbf{x} , the black-box model with input \mathbf{x} produces the corresponding output in the form of the CHF $H(t|\mathbf{x})$ or the hazard function $h(t|\mathbf{x})$. The basic idea behind the explanation algorithm SurvLIME is to approximate the output of the black-box model in the form of the CHF by means of the CHF produced by the Cox model for the same input example \mathbf{x} . The idea stems from the fact that the function $\psi(\mathbf{x}, \mathbf{b})$ in the Cox model (see (6)) is linear and does not depend on time t . The linearity means that coefficients b_i of the covariates in $\psi(\mathbf{x}, \mathbf{b})$ can be regarded as quantitative impacts on the prediction. Hence, the largest coefficients indicate the corresponding importance features. The independence of $\psi(\mathbf{x}, \mathbf{b})$ on time t means that time and covariates can be considered separately, and the optimization problem minimizing the difference between CHFs is significantly simplifies.

In order to find the important features of \mathbf{x} , we have to compute some optimal values of elements of vector \mathbf{b} (see the previous section) of the Cox model such that $H(t|\mathbf{x})$ would be as close as possible to the Cox CHF denoted as $H_{\text{Cox}}(t|\mathbf{x}, \mathbf{b})$. However, the use of a single point may lead to incorrect results. Therefore, we generate a lot of nearest points \mathbf{x}_k in a local area around \mathbf{x} . For every generated point \mathbf{x}_k , the CHF $H(t|\mathbf{x}_k)$ of the black-box model can be computed as a prediction provided by the survival black-box model, and the Cox CHF $H_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$ can be written as a function of unknown vector \mathbf{b} . Now the optimal values of \mathbf{b} can be computed by minimizing the average distance between every pair of CHFs $H(t|\mathbf{x}_k)$ and $H_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$ over all generated points \mathbf{x}_k . Every distance between CHFs has a weight w_k which depends on the distance between \mathbf{x}_k and \mathbf{x} . Smaller distances between \mathbf{x}_k and \mathbf{x} produce larger weights of distances between CHFs. If explanation methods like LIME deal with point-valued predictions of the example processing through the black-box model, then the proposed method is based on functions characterizing every point from D or new points \mathbf{x} . Fig. 1 illustrates the explanation algorithm. It can be seen from Fig. 1 that a set of examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ are fed to the black-box survival model, which produces a set of CHFs $\{H(t|\mathbf{x}_1), \dots, H(t|\mathbf{x}_N)\}$. Simultaneously, we write CHFs $H_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$, $k = 1, \dots, N$, as functions of coefficients \mathbf{b} for all generated examples. The weighted average distance between the CHFs of the Cox model and the black-box survival model allows us to construct an optimization problem and to compute the optimal vector \mathbf{b} by solving the optimization problem.

It should be noted that the above description is only a sketch where every step is a separate task. All steps will be considered in detail below.

6 Minimization of distances between functions

It has been mentioned that the main peculiarity of machine learning survival models is that the output of models is a function (the CHF or the survival function). Therefore, in order to approximate the output of the black-box model by means of the CHF produced by the Cox model at the input example \mathbf{x} , we have to generate many points \mathbf{x}_k in a local area around \mathbf{x} and to consider the mean distance between the CHFs for generated points \mathbf{x}_k , $k = 1, \dots, N$, and the Cox model CHF for point \mathbf{x} . Before deriving the mean distance and its minimizing, we introduce some notations and conditions.

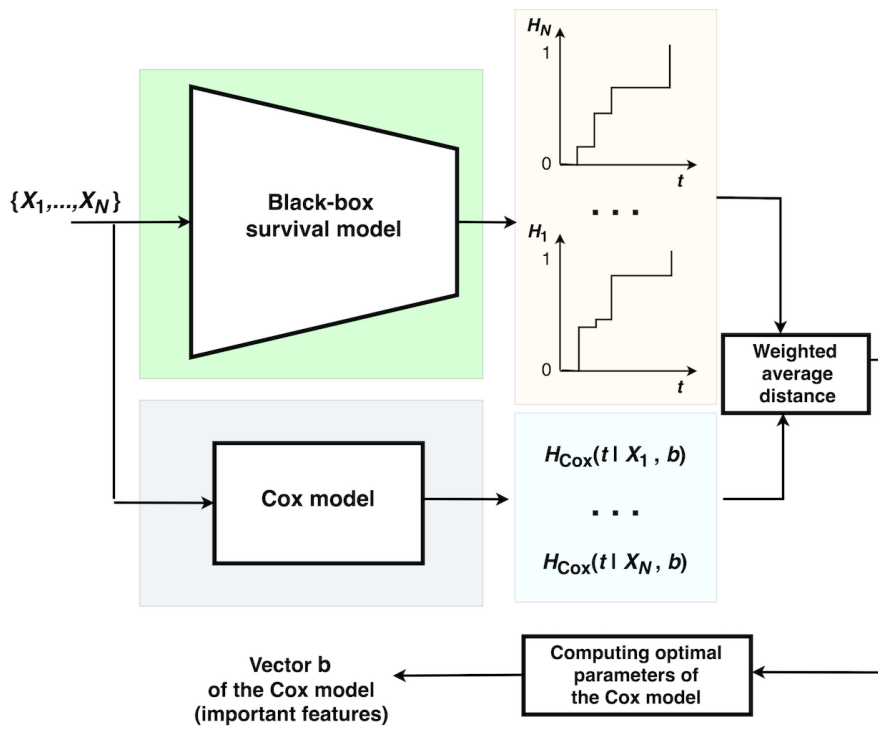


Figure 1: A schematic illustration of the explanation algorithm

Let $t_0 < t_1 < \dots < t_m$ be the distinct times to event of interest, for example, times to deaths from the set $\{T_1, \dots, T_n\}$, where $t_0 = \min_{k=1, \dots, n} T_k$ and $t_m = \max_{k=1, \dots, n} T_k$. The black-box model maps the feature vectors $\mathbf{x} \in \mathbb{R}^d$ into piecewise constant CHF's $H(t|\mathbf{x})$ having the following properties:

1. $H(t|\mathbf{x}) \geq 0$ for all t
2. $\max_t H(t|\mathbf{x}) < \infty$
3. $\int_0^\infty H(t|\mathbf{x}) dt \rightarrow \infty$

Let us introduce the time $T = t_m + \gamma$ in order to restrict the integral of $H(t|\mathbf{x})$, where γ is a very small positive number. Let $\Omega = [0, T]$. Then we can write

$$\int_{\Omega} H(t|\mathbf{x}) dt < \infty. \quad (11)$$

Since the CHF $H(t|\mathbf{x})$ is piecewise constant, then it can be written in a special form. Let us divide the set Ω into $m + 1$ subsets $\Omega_0, \dots, \Omega_m$ such that

1. $\Omega = \cup_{j=0, \dots, m} \Omega_j$
2. $\Omega_m = [t_m, T]$, $\Omega_j = [t_j, t_{j+1})$, $\forall j \in \{0, \dots, m - 1\}$
3. $\Omega_j \cap \Omega_k = \emptyset$ for $\forall j \neq k$

Let us introduce the indicator functions

$$\chi_{\Omega_j}(t) = \begin{cases} 1, & t \in \Omega_j, \\ 0, & t \notin \Omega_j. \end{cases} \quad (12)$$

Hence, the CHF $H(t|\mathbf{x})$ can be expressed through the indicator functions as follows:

$$H(t|\mathbf{x}) = \sum_{j=0}^m H_j(\mathbf{x}) \cdot \chi_{\Omega_j}(t) \quad (13)$$

under additional condition $H_j(\mathbf{x}) \geq \varepsilon > 0$, where ε is a small positive number. Here $H_j(\mathbf{x})$ is a part of the CHF in interval Ω_j . It is important that $H_j(\mathbf{x})$ does not depend on t and is constant in interval Ω_j . The last condition will be necessary below in order to deal with logarithms of the CHF's.

Let g be a monotone function. Then there holds

$$g(H(t|\mathbf{x})) = \sum_{j=0}^m g(H_j(\mathbf{x})) \chi_{\Omega_j}(t). \quad (14)$$

By using the above representation of the CHF and integrating it over Ω , we get

$$\begin{aligned} \int_{\Omega} H(t|\mathbf{x}) dt &= \int_{\Omega} \left[\sum_{j=0}^m H_j(\mathbf{x}) \chi_{\Omega_j}(t) \right] dt \\ &= \sum_{j=0}^m H_j(\mathbf{x}) \left[\int_{\Omega} \chi_{\Omega_j}(t) dt \right] = \sum_{j=0}^m H_j(\mathbf{x}) (t_{j+1} - t_j). \end{aligned} \quad (15)$$

The same expressions can be written for the Cox CHF:

$$H_{\text{Cox}}(t|\mathbf{x}, \mathbf{b}) = H_0(t) \exp(\mathbf{b}^T \mathbf{x}) = \sum_{j=0}^m [H_{0j} \exp(\mathbf{b}^T \mathbf{x})] \chi_{\Omega_j}(t), \quad H_{0j} \geq \varepsilon. \quad (16)$$

It should be noted that the distance between two CHFs can be replaced with the distance between two logarithms of the corresponding CHFs for the optimization problem. The introduced condition $H_j(\mathbf{x}) \geq \varepsilon > 0$ allows to use logarithms. Therefore, in order to get the convex optimization problem for finding optimal values of \mathbf{b} , we consider logarithms of the CHFs. It is important to point out that the difference of logarithms of the CHFs is not equal to the difference between CHFs themselves. However, we make this replacement to simplify the optimization problem for computing important features.

Let $\phi(t|\mathbf{x}_k)$ and $\phi_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$ be logarithms of $H(t|\mathbf{x}_k)$ and $H_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$. Here \mathbf{x}_k is a generated point. The difference between functions $\phi(t|\mathbf{x}_k)$ and $\phi_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$ can be written as follows:

$$\begin{aligned} & \phi(t|\mathbf{x}_k) - \phi_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b}) \\ &= \sum_{j=0}^m (\ln H_j(\mathbf{x}_k)) \chi_{\Omega_j}(t) - \sum_{j=0}^m (\ln(H_{0j} \exp(\mathbf{b}^T \mathbf{x}_k))) \chi_{\Omega_j}(t) \\ &= \sum_{j=0}^m (\ln H_j(\mathbf{x}_k) - \ln H_{0j} - \mathbf{b}^T \mathbf{x}_k) \chi_{\Omega_j}(t). \end{aligned} \quad (17)$$

Let us consider the distance between functions $\phi(t|\mathbf{x}_k)$ and $\phi_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})$ in metric L_2 :

$$\begin{aligned} D_{2,k}(\phi, \phi_{\text{Cox}}) &= \|\phi(t|\mathbf{x}_k) - \phi_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})\|_2^2 \\ &= \int_{\Omega} |\phi(t|\mathbf{x}_k) - \phi_{\text{Cox}}(t|\mathbf{x}_k, \mathbf{b})|^2 dt \\ &= \sum_{j=0}^m (\ln H_j(\mathbf{x}_k) - \ln H_{0j} - \mathbf{b}^T \mathbf{x}_k)^2 (t_{j+1} - t_j). \end{aligned} \quad (18)$$

The function $\ln H_j(\mathbf{x}_k) - \ln H_{0j} - \mathbf{b}^T \mathbf{x}_k$ is linear with \mathbf{b} and, therefore, convex. Moreover, it is easy to prove by taking the second derivative over b_j that the function $(A - \mathbf{b}^T \mathbf{x}_k)^2$ is also convex, where $A = \ln H_j(\mathbf{x}_k) - \ln H_{0j}$. Since the term $(t_{j+1} - t_j)$ is positive, then the function $D_{2,k}(\phi, \phi_{\text{Cox}})$ is also convex as a linear combination of convex functions with positive coefficients.

According to the explanation algorithm, we have to consider many points \mathbf{x}_k generated in a local area around \mathbf{x} and to minimize the objective function $\sum_{k=1}^N w_k D_{2,k}(\phi, \phi_{\text{Cox}})$, which takes into account all these generated examples and the corresponding weights of the examples. It is obvious that this function is also convex with respect to \mathbf{b} . The weight can be assigned to point \mathbf{x}_k as a decreasing function of the distance between \mathbf{x} and \mathbf{x}_k , for example, $w_k = K(\mathbf{x}, \mathbf{x}_k)$, where $K(\cdot, \cdot)$ is a kernel. In numerical experiments, we use the function defined in (24).

Finally, we write the following convex optimization problem:

$$\min_{\mathbf{b}} \sum_{k=1}^N w_k \sum_{j=0}^m (\ln H_j(\mathbf{x}_k) - \ln H_{0j} - \mathbf{b}^T \mathbf{x}_k)^2 (t_{j+1} - t_j). \quad (19)$$

One of the difficulties of solving the above problem is that the difference between functions $H(t|\mathbf{x})$ and $H_{\text{Cox}}(t|\mathbf{x}, \mathbf{b})$ may be significantly different from the distance between their logarithms. Therefore,

in order to take into account this fact, it is proposed to introduce weights which “straighten” functions $\phi(t|\mathbf{x})$ and $\phi_{\text{Cox}}(t|\mathbf{x}, \mathbf{b})$. These weights are defined as

$$v(t|\mathbf{x}) = \frac{H(t|\mathbf{x})}{\phi(t|\mathbf{x})} = \frac{H(t|\mathbf{x})}{\ln(H(t|\mathbf{x}))}. \quad (20)$$

Taking into account these weights and their representation $v_{kj} = H_j(\mathbf{x}_k)/\ln(H_j(\mathbf{x}_k))$, the optimization problem can be rewritten as

$$\min_{\mathbf{b}} \sum_{k=1}^N w_k \sum_{j=0}^m v_{kj}^2 (\ln H_j(\mathbf{x}_k) - \ln H_{0j} - \mathbf{b}^T \mathbf{x}_k)^2 (t_{j+1} - t_j). \quad (21)$$

The problem can be solved by many well-known methods of the convex programming. Finally, we write the following scheme of Algorithm 1.

Algorithm 1 The algorithm for computing vector \mathbf{b} for point \mathbf{x}

Require: Training set D ; point of interest \mathbf{x} ; the number of generated points N ; the black-box survival model for explaining $f(\mathbf{x})$

Ensure: Vector \mathbf{b} of important features

- 1: Compute the baseline cumulative hazard function $H_0(t)$ of the approximating Cox model on dataset D by using the Nelson–Aalen estimator
 - 2: Generate $N - 1$ random nearest points \mathbf{x}_k in a local area around \mathbf{x} , point \mathbf{x} is the N -th point
 - 3: Get the prediction of the cumulative hazard function $H(t|\mathbf{x}_k)$ by using the black-box survival model (the function f)
 - 4: Compute weights $w_k = K(\mathbf{x}, \mathbf{x}_k)$ of perturbed points, $k = 1, \dots, N$
 - 5: Compute weights $v_{kj} = H_j(\mathbf{x}_k)/\ln(H_j(\mathbf{x}_k))$, $k = 1, \dots, N$, $j = 0, \dots, m$
 - 6: Find vector \mathbf{b} by solving the convex optimization problem (21)
-

7 Numerical experiments with synthetic data

In order to study the proposed explanation algorithm, we generate random survival times to events by using the Cox model estimates. This generation allows us to compare initial data for generating every points and results of SurvLIME.

7.1 Generation of random covariates, survival times and perturbations

Two clusters of covariates $\mathbf{x} \in \mathbb{R}^d$ are randomly generated such that points of every cluster are generated from the uniform distribution in a sphere. The covariate vectors are generated in the d -sphere with some predefined radius R . The center p of the d -sphere and its radius R are parameters of experiments. There are several methods for the uniform sampling of points \mathbf{x} in the d -sphere with the unit radius $R = 1$, for example, [5, 20]. Then every generated point is multiplied by R . The following parameters for points of every cluster are used:

1. cluster 0: center $p_0 = (0, 0, 0, 0, 0)$; radius $R = 8$; number of points $N = 1000$;
2. cluster 1: center $p_1 = (4, -8, 2, 4, 2)$; radius $R = 8$; number of points $N = 1000$.

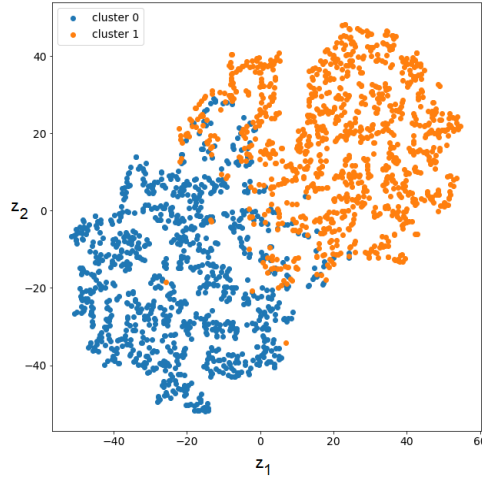


Figure 2: Two clusters of generated covariates depicted by using the t-SNE method

The parameters of clusters are chosen in order to get some intersecting area containing points from both the clusters, in particular, the radius is determined from the centers as follows:

$$R = \left\lceil \frac{\|p_0 - p_1\|_2}{2} \right\rceil + 2. \quad (22)$$

The clusters after using the well-known t-SNE algorithm are depicted in Fig. 2.

In order to generate random survival times by using the Cox model estimates, we apply results obtained by Bender et al. [7] for survival time data for the Cox model with Weibull distributed survival times. The Weibull distribution with the scale λ and shape v parameters is used to generate appropriate survival times for simulation studies because this distribution shares the assumption of proportional hazards with the Cox regression model [7]. Taking into account the parameters λ and v of the Weibull distribution, we use the following expression for generated survival times [7]:

$$T = \left(\frac{-\ln(U)}{\lambda \exp(\mathbf{b}^T \mathbf{x})} \right)^{1/v}, \quad (23)$$

where U is the random variable uniformly distributed in interval $[0, 1]$.

Parameters of the generation for every cluster are

1. cluster 0: $\lambda = 10^{-5}$, $v = 2$, $\mathbf{b}^T = (10^{-6}, 0.1, -0.15, 10^{-6}, 10^{-6})$;
2. cluster 1: $\lambda = 10^{-5}$, $v = 2$, $\mathbf{b}^T = (10^{-6}, -0.15, 10^{-6}, 10^{-6}, -0.1)$.

It can be seen from the above that every vector \mathbf{b} has three almost zero-valued elements and two “large” elements which will correspond to important features. Generated values T_i are restricted by the condition: if $T_i > 2000$, then T_i is replaced with value 2000. The event indicator δ_i is generated from the binomial distribution with probabilities $\Pr\{\delta_i = 1\} = 0.9$, $\Pr\{\delta_i = 0\} = 0.1$.

Perturbation is one of the steps of the algorithm. According to it, we generate N nearest points \mathbf{x}_k in a local area around \mathbf{x} . These points are uniformly generated in the d -sphere with some predefined

radius $r = 0.5$ and the center at point \mathbf{x} . In numerical experiments, $N = 1000$. Weights to every point are assigned as follows:

$$w_k = 1 - \sqrt{\frac{\|\mathbf{x} - \mathbf{x}_k\|_2}{r}}. \quad (24)$$

7.2 Black-box models and approximation measures

As black-box models, we use the Cox model and the RSF model [26]. The RSF consists of 250 decision survival trees. The approximating Cox model has the baseline CHF $H_0(t)$ constructed on generated training data using the Nelson–Aalen estimator. The Cox model is used in order to check whether the selected important features explaining the CHF $H(t|\mathbf{x})$ at point \mathbf{x} coincide with the corresponding features accepted in the Cox model for generating training set. It should be noted that the Cox model as well as the RSF are viewed as black-box models whose predictions (CHF's or survival functions) are explained. To study how different cases impact on the quality of the approximation, we use the following two measures for the Cox model:

$$RMSE_{\text{model}} = \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{b}_i^{\text{model}} - \mathbf{b}_i^{\text{expl}}\|_2}, \quad (25)$$

$$RMSE_{\text{true}} = \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{b}_i^{\text{true}} - \mathbf{b}_i^{\text{expl}}\|_2}, \quad (26)$$

where $\mathbf{b}_i^{\text{model}}$ are coefficients of the Cox model which is used as the black-box model; $\mathbf{b}_i^{\text{true}}$ are coefficients used for training data generation (see (23)); $\mathbf{b}_i^{\text{expl}}$ are explaining coefficients obtained by using the proposed algorithm.

The first measure characterizes how the obtained important features coincide with the corresponding features obtained by using the Cox model as the black-box model. The second measure considers how the obtained important feature coincide with the features used for generating the random times to events. Every measure is calculated by taking randomly n points \mathbf{x} from the testing set and compute the corresponding coefficients.

In order to investigate the quality of explanation when the black-box model is the RSF, we use another measure:

$$RMSE_{\text{approx}} = \sqrt{\frac{1}{n} \sum_{i=1}^n \sum_{j \in J} \left(H(t_j|\mathbf{x}_i) - H_{\text{Cox}}(t_j|\mathbf{x}_i, \mathbf{b}_i^{\text{expl}}) \right)^2}, \quad (27)$$

where J is a set of time indices for computing the measure.

This measure considers how the obtained Cox model approximation $H_{\text{Cox}}(t_j|\mathbf{x}_i, \mathbf{b}_i^{\text{expl}})$ is close to the RSF output $H(t_j|\mathbf{x}_i)$. We cannot estimate the proximity of important features explaining the model because we do not have the corresponding features for the RSF. A comparison of the important features with the generated ones is also incorrect because we explain the model (the RSF) output ($H(t|\mathbf{x}_i)$), but not training data. Therefore, we use the above measure to estimate the proximity of two CHF's: the Cox model approximation and the RSF output.

7.3 Experiment 1

To evaluate the algorithm, 900 examples are randomly selected from every cluster for training and 100 examples are for testing. Three cases of training and testing the black-box Cox and RSF models are studied:

1. cluster 0 for training and for testing;
2. cluster 1 for training and for testing;
3. clusters 0 and 1 are jointly used for training and separately for testing.

The testing phase includes:

- Computing the explanation vector \mathbf{b}^{expl} for every point from the testing set.
- Depicting the best, mean and worst approximations in accordance with Euclidean distance between vectors \mathbf{b}^{expl} and $\mathbf{b}^{\text{model}}$ (for the Cox model) and with Euclidean distance between $H(t_j|\mathbf{x}_i)$ and $H_{\text{Cox}}(t_j|\mathbf{x}_i, \mathbf{b}_i^{\text{expl}})$ (for the RSF). In order to get these approximations, points with the best, mean and worst approximations are selected among all testing points.
- Computing measures $RMSE_{\text{model}}$ and $RMSE_{\text{true}}$ for the Cox model and $RMSE_{\text{approx}}$ for the RSF over all points of the testing set.

The three cases (best (pictures in the first row), mean (pictures in the second row) and worst (pictures in the third row)) of approximations for the black-box Cox model under condition that cluster 0 is used for training and testing are depicted in Fig. 3. Left pictures show values of important features \mathbf{b}^{expl} , $\mathbf{b}^{\text{model}}$ and \mathbf{b}^{true} . It can be seen from these pictures that all experiments show very clear coincidence of important features for all models. Right pictures in Fig. 3 show survival functions obtained from the black-box Cox model and from the Cox approximation. It follows from the pictures that the approximation is perfect even for the worst case. Similar results can be seen from Fig. 4, where training and testing examples are taken from cluster 1.

Figs. 5 and 6 illustrate different results corresponding to cases when training examples are taken from cluster 0 and cluster 1. One can see that the approximation of survival functions is perfect, but important features obtained by SurvLIME do not coincide with the features used for generating random times in accordance with the Cox model. This fact can be explained in the following way. We try to explain the CHF obtained by using the black-box model (the Cox model in the considered case). But the black-box Cox model is trained on all examples from two different clusters. We have a mix of data from two clusters with different parameters. Therefore, this model itself provides results different from the generation model. At the same time, one can observe from Figs. 5 and 6 that the explaining important features coincide with the features which are obtained from the black-box Cox model. These results are interesting. They show that we explain CHFs of the black-box model, but not CHFs of training data.

The approximation accuracy measures for four cases are given in Table 1. In fact, the table repeats the results shown in Figs. 3-6.

In the same way, we study the black-box RSF model by using three cases of its training and testing. The results are shown in Fig. 7 where the first, second, third, and fourth rows correspond to the four cases of experiments: cluster 0 for training and for testing; cluster 1 for training and for testing; clusters 0 and 1 are jointly used for training and separately for testing. Every row contains

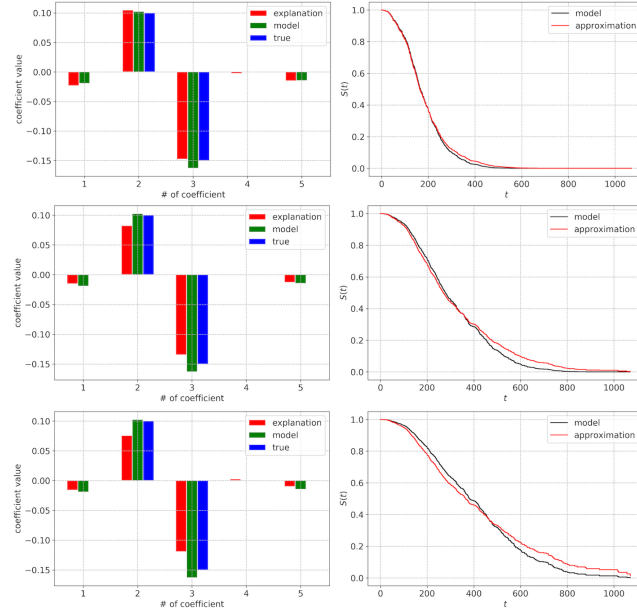


Figure 3: The best, mean and worst approximations for the Cox model (training and testing sets from cluster 0)

Table 1: Approximation accuracy measures for four cases of using the black-box Cox model

Clusters for training	Cluster for testing	$RMSE_{\text{model}}$	$RMSE_{\text{true}}$	
			0	1
0	0	0.016	0.014	
1	1	0.038		0.048
$\{0, 1\}$	0	0.023	0.089	
$\{0, 1\}$	1	0.014		0.065

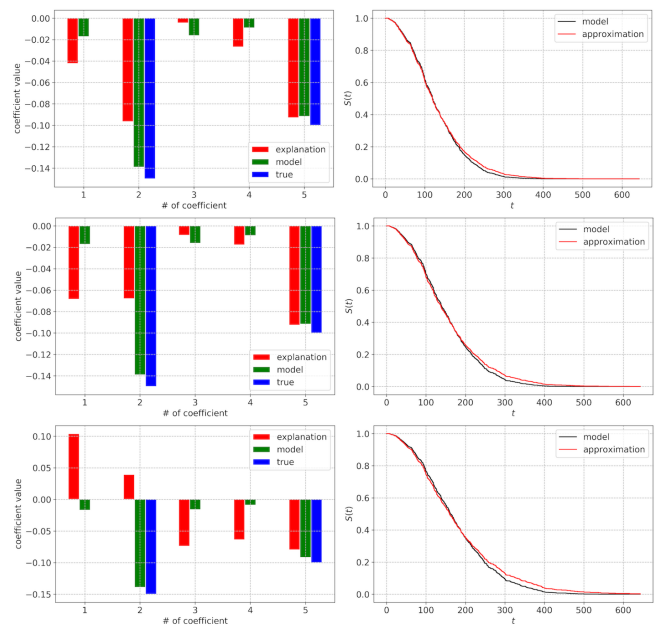


Figure 4: The best, mean and worst approximations for the Cox model (training and testing sets from cluster 1)

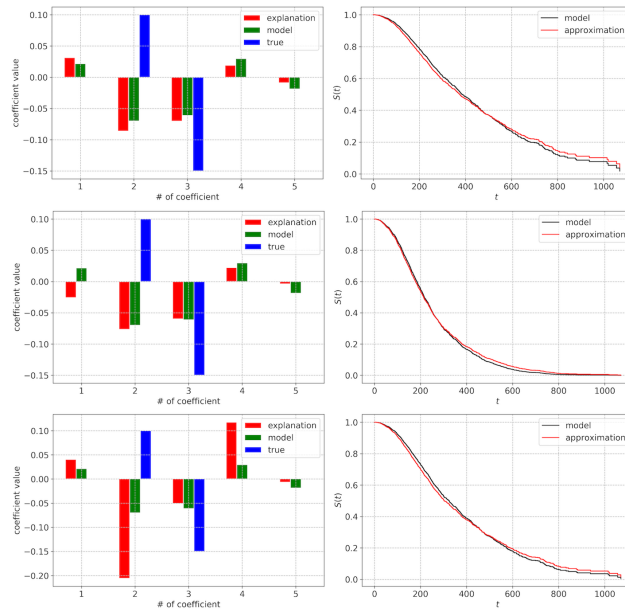


Figure 5: The best, mean and worst approximations for the Cox model (the training set from clusters 0,1 and the testing set from cluster 0)

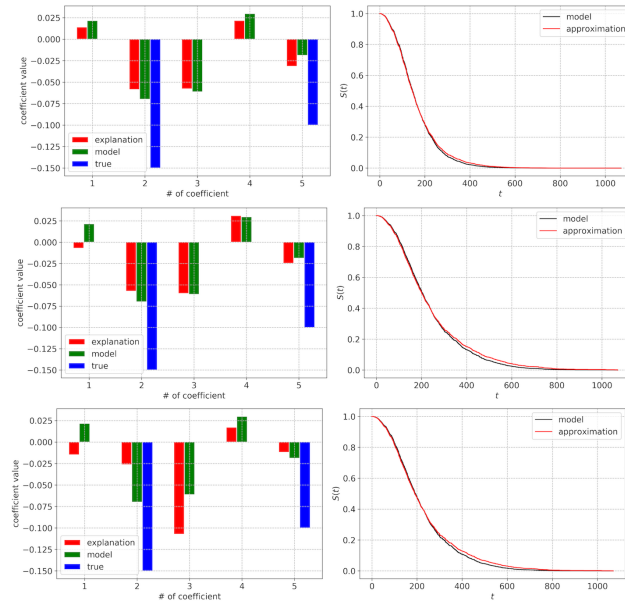


Figure 6: The best, mean and worst approximations for the Cox model (the training set from clusters 0,1 and the testing set from cluster 1)

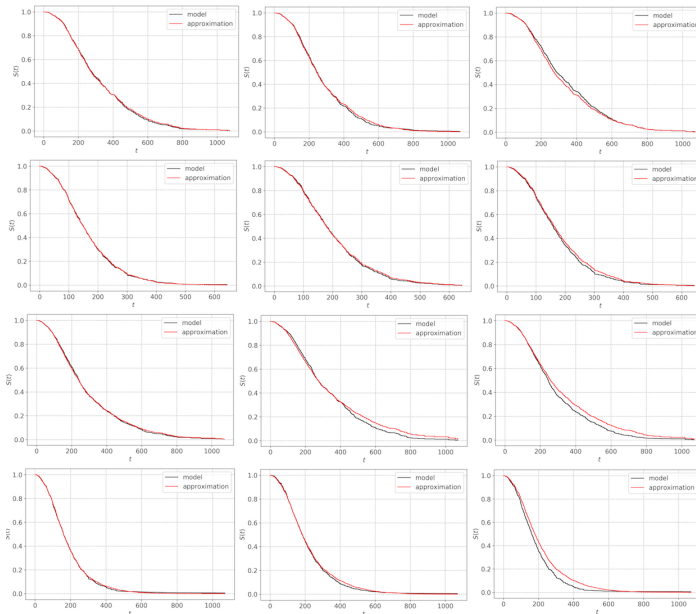


Figure 7: The best, mean and worst approximations for the RSF model

pairs of survival functions corresponding to the best, mean and worst approximations. The important features are not shown in Fig. 7 because they cannot be compared with the features of the generative Cox model. Moreover, the RSF does not provide the important features like the Cox model. One can again see from Fig. 7 that SurvLIME illustrates the perfect approximation of the RSF output by the Cox model.

7.4 Experiment 2

In the second experiment, our aim is to study how SurvLIME depends on the number of training examples. We use only the Cox model for explaining which is trained on 100, 200, 300, 400, 500 examples and tested on 100 examples from cluster 0. We study how the difference between $\mathbf{b}^{\text{model}}$ and \mathbf{b}^{true} depends of the sample size. Results of experiments are shown in Fig. 8 where rows correspond to 100, 200, 300, 400, 500 training examples, respectively, left pictures illustrate relationships between important features, right pictures show survival functions obtained from the black-box Cox model and from the Cox approximation (see similar pictures in Figs. 3-6). It is interesting to observe in Fig. 8 how the vectors \mathbf{b}^{expl} , $\mathbf{b}^{\text{model}}$ and \mathbf{b}^{true} approach each other with the number of training examples.

The measures $RMSE_{\text{model}}$ and $RMSE_{\text{true}}$ as functions of the sample size are provided in Table 2. It can be seen from Table 2 a tendency of the measures to be reduced with increase of the sample size for training. The C-index as a measure of the black-box model quality is also given in Table 2.

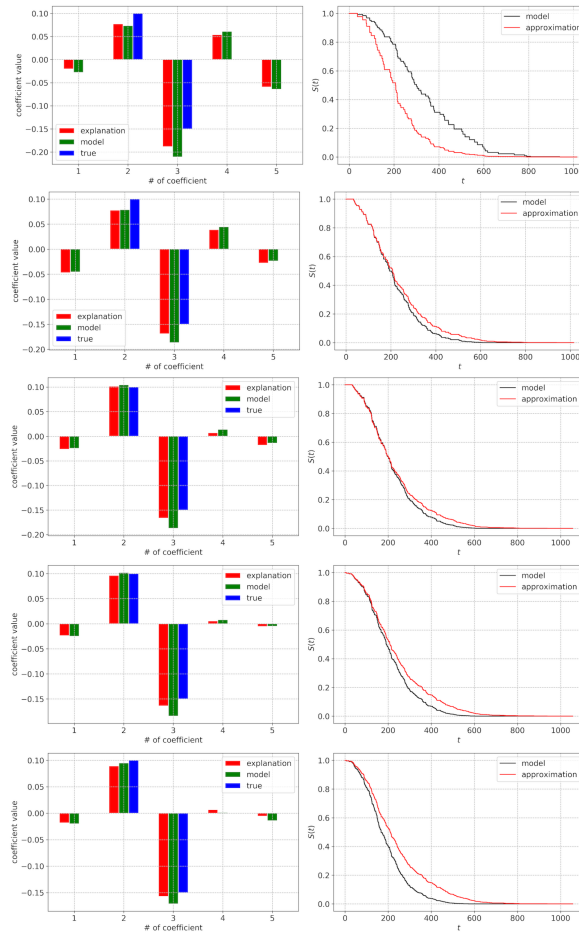


Figure 8: Comparison of important features (left pictures) and survival functions (right pictures) for the black-box Cox model by 100, 200, 300, 400, 500 training examples

Table 2: Approximation accuracy measures for four cases of using the black-box Cox model

	C-index	$RMSE_{\text{model}}$	$RMSE_{\text{true}}$
100	0.777	0.027	0.042
200	0.785	0.018	0.031
300	0.787	0.021	0.015
400	0.838	0.019	0.014
500	0.844	0.017	0.015

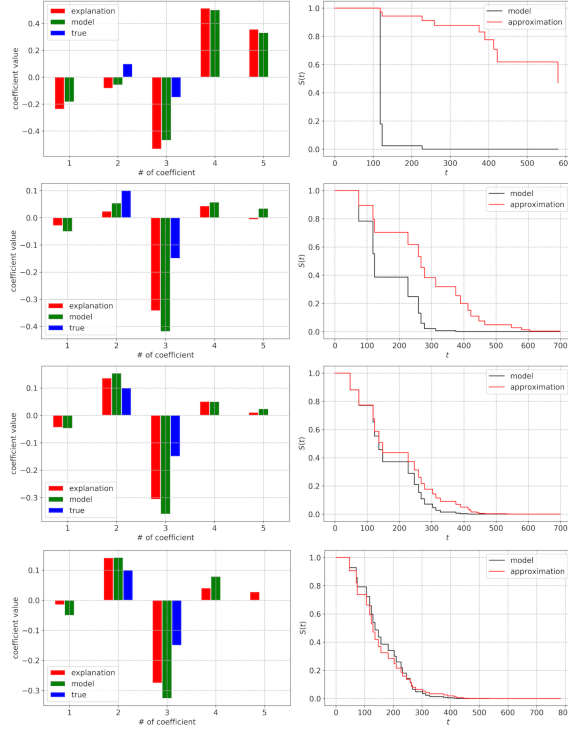


Figure 9: Comparison of important features (left pictures) and survival functions (right pictures) for the black-box Cox model by 10, 20, 30, 40 training examples

7.5 Experiment 3

Another interesting question for studying is how SurvLIME behaves by a very small amount of training data. We use the Cox model and the RSF as black-box models and train them on 10, 20, 30, 40 examples from cluster 0. The models are tested on 10 examples from cluster 0. Results of experiments for the Cox model are shown in Fig. 9 where rows correspond to 10, 20, 30, 40 training examples, respectively, left pictures illustrate relationships between important features, right pictures show survival functions obtained from the black-box Cox model and from the Cox approximation (see similar pictures in Figs. 8). It is interesting to point out from Fig. 9 that important features obtained by means of SurvLIME differ from the parameters of the Cox model for generating random examples in the case of 10 training examples. However, they almost coincide with features obtained by the black-box model. Moreover, it follows from Fig. 9 that important features of SurvLIME quickly converge to parameters of the Cox model for generating random examples.

The C-index, measures $RMSE_{\text{model}}$ and $RMSE_{\text{true}}$ as functions of the sample size are provided in Table 3. It can be seen from Table 3 that $RMSE_{\text{model}}$ and $RMSE_{\text{true}}$ are strongly reduced with increase of the sample size.

Results of experiments for the RSF are shown in Fig. 10 where rows correspond to 10, 20, 30, 40 training examples, respectively, left pictures illustrate important features, right pictures show survival

Table 3: Approximation accuracy measures for four cases of using the black-box Cox model by the small amount of data

	C-index	$RMSE_{\text{model}}$	$RMSE_{\text{true}}$
10	0.444	0.231	0.293
20	0.489	0.194	0.197
30	0.622	0.054	0.082
40	0.733	0.047	0.053

Table 4: A brief introduction about datasets

Data set	Abbreviation	R Package	d	d^*	m
Chronic Myelogenous Leukemia Survival	CML	multcomp	5	9	507
NCCTG Lung Cancer	LUNG	survival	8	11	228
Primary Biliary Cirrhosis	PBC	survival	17	22	418
Stanford Heart Transplant	Stanford2	survival	2	2	185
Trial Of Usrodeoxycholic Acid	UDCA	survival	4	4	170
Veterans' Administration Lung Cancer Study	Veteran	survival	6	9	137

functions obtained from the black-box RSF and from the Cox approximation (see similar pictures in Fig. 9). It can be seen from Fig. 10 that the difference between survival functions is reduced with increase the training set. However, in contrast to the black-box Cox model (see Fig. 9), important features are very unstable. They are different for every training sets. A reason of this important feature behavior is that they explain the RSF outputs which are significantly changed by the so small numbers of training examples.

8 Numerical experiments with real data

In order to illustrate SurvLIME, we test it on several well-known real datasets. A short introduction of the benchmark datasets are given in Table 4 that shows sources of the datasets (R packages), the number of features d for the corresponding dataset, the number of training instances m and the number of extended features d^* using hot-coding for categorical data.

Figs. 11-12 illustrate numerical results for the CML dataset. Three cases of approximation are considered: best (pictures in the first row), mean (pictures in the second row) and worst (pictures in the third row). These cases are similar to cases studied for synthetic data. The cases are studied for the black-box Cox model (Fig. 11) and the black-box RSF (Fig. 12). Again left pictures in figures show values of important features $\mathbf{b}^{\text{model}}$ and \mathbf{b}^{true} for the Cox model and \mathbf{b}^{true} for the RSF, right pictures illustrate the approximated survival function and the survival function obtained by the explained model.

Figs. 13-17 show numerical results for other datasets. Since most results are very similar to the same results obtained for the CML datasets, then we provide only the case of the mean approximation for every dataset in order to reduce the number of similar pictures. Moreover, we do not show important features explaining RSFs because, in contrast to the Cox model, they cannot be compared with true features. Every figure consists of three pictures: the first one illustrates the explanation important features and important features obtained from training the Cox model; the second picture shows two survival functions for the Cox model; the third picture shows two survival functions for the

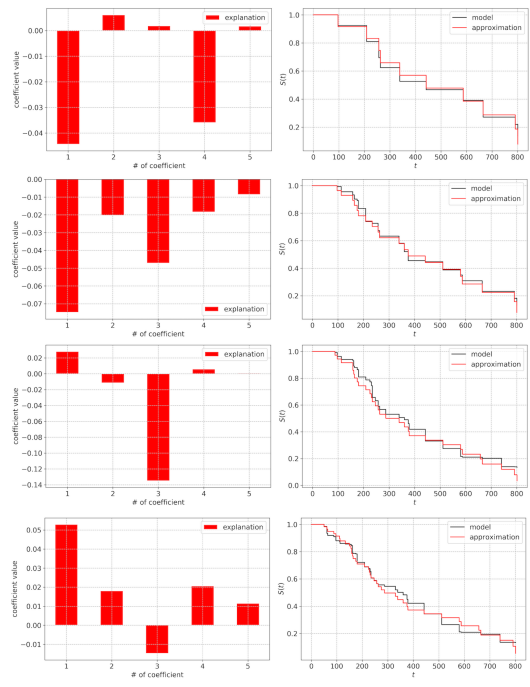


Figure 10: Important features (left pictures) and survival functions (right pictures) for the black-box RSF by 10, 20, 30, 40 training examples

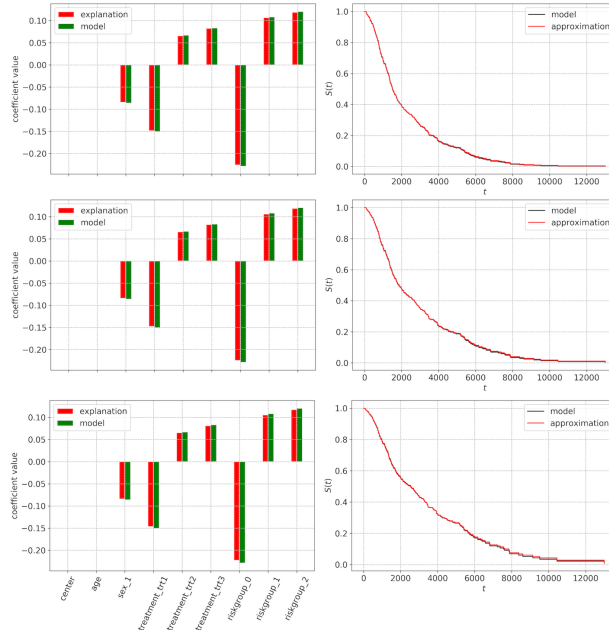


Figure 11: The best, mean and worst approximations for the Cox model trained on the CML dataset

RSF.

If the Cox model is used for training, then one can see an explicit coincidence of explaining important features and the features obtained from the trained Cox model for all datasets. This again follows from the fact that SurvLIME does not explain a dataset. It explains the model, in particular, the Cox model. We also can observe that the approximated survival function is very close to the survival function obtained by the explained Cox model in all cases. The same cannot be said about the RSF. For example, one can see from Fig. 12 that the important features obtained by explaining the RSF mainly do not coincide. The reason is a difference of results provided by the Cox model and the RSF. At the same time, it can be seen from Figs. 12-17 that the survival functions obtained from the RSF and the approximating Cox model are close to each other for many models. This implies that SurvLIME provides correct results.

9 Conclusion

A new explanation method called SurvLIME which can be regarded as a modification of the well-known method LIME for survival data has been presented in the paper. The main idea behind the method is to approximate a survival machine learning model at a point by the Cox proportional hazards model which assumes a linear combination of the example covariates. This assumption allows us to determine the important features explaining the survival model.

In contrast to LIME and other explanation methods, SurvLIME deals with survival data. It is not complex from computational point of view. Indeed, we have derived a simple convex unconstrained optimization problem whose solution does not meet any difficulties. Moreover, many numerical exper-

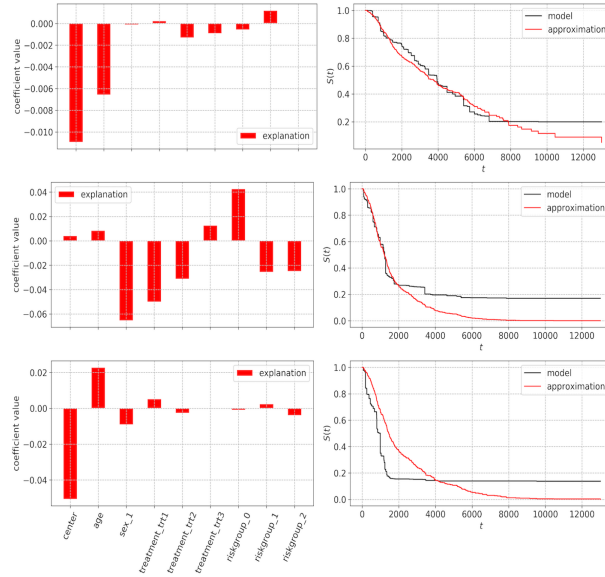


Figure 12: The best, mean and worst approximations for the RSF trained on the CML dataset

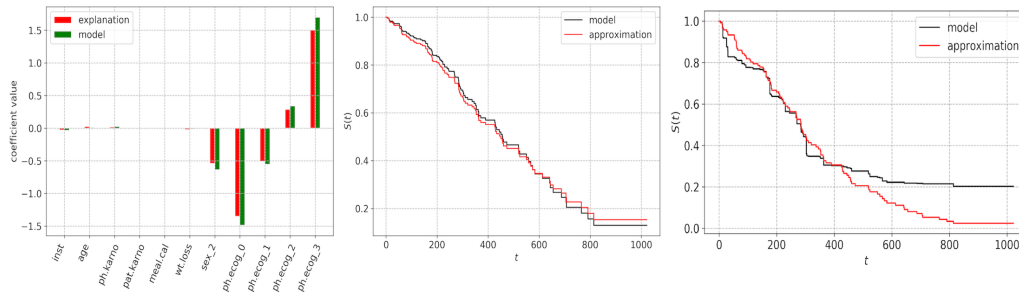


Figure 13: The mean approximation for the Cox model (the first and the second picture) and the RSF (the third picture) trained on the LUNG dataset

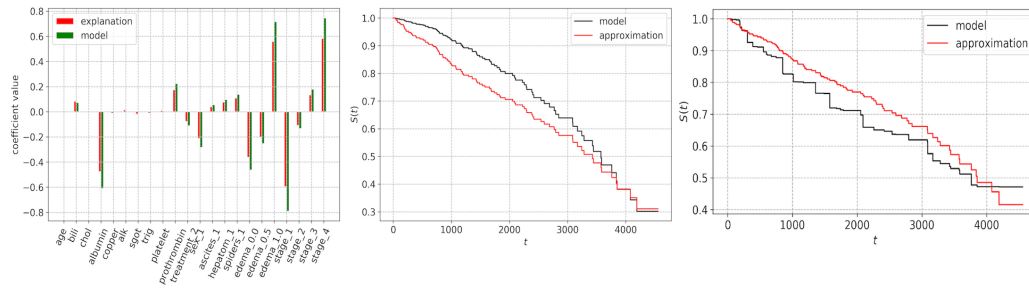


Figure 14: The mean approximation for the Cox model (the first and the second picture) and the RSF (the third picture) trained on the PBC dataset

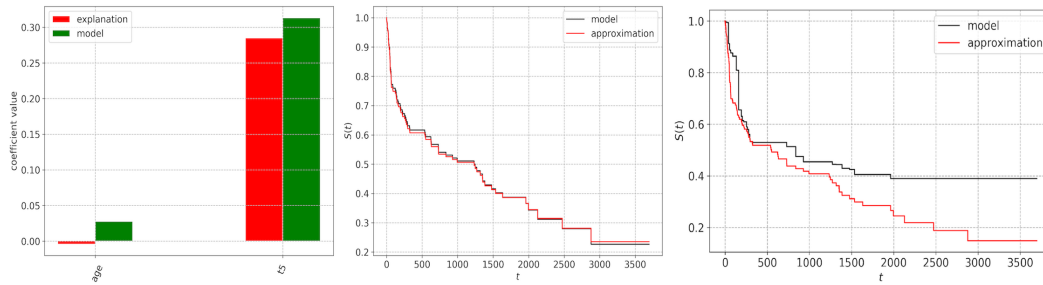


Figure 15: The mean approximation for the Cox model (the first and the second picture) and the RSF (the third picture) trained on the Stanford2 dataset

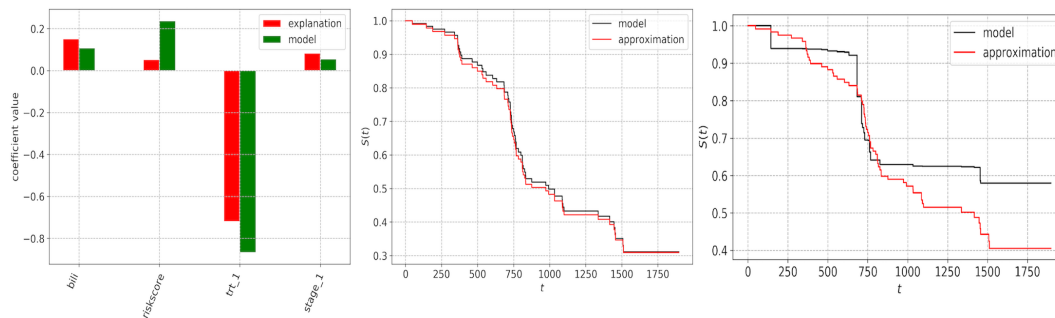


Figure 16: The mean approximation for the Cox model (the first and the second picture) and the RSF (the third picture) trained on the UDCA dataset

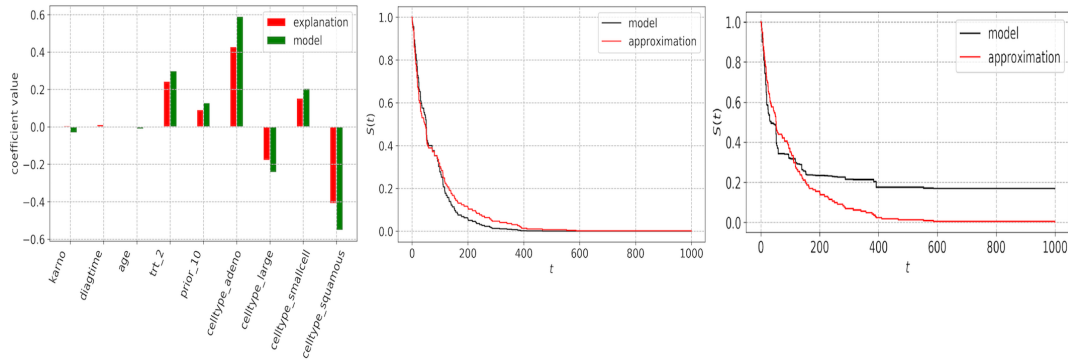


Figure 17: The mean approximation for the Cox model (the first and the second picture) and the RSF (the third picture) trained on the Veteran dataset

iments with synthetic and real datasets have clearly illustrated accuracy and correctness of SurvLIME. It has coped even with problems characterizing by small datasets.

The main advantage of the method is that it opens a door for developing many explanation methods taking into account censored data. In particular, an interesting problem is to develop a method explaining the survival models with time-dependent covariates. This is a direction for further research. Only the quadratic norm has been considered to estimate the distance between two CHF's and to construct the corresponding optimization problem. However, there are other distance metrics which are also interesting with respect to constructing new explanation methods. This is another direction for further research. An open and very interesting direction is also the counterfactual explanation with censored data. It could be a perspective extension of SurvLIME. It should be noted that SurvLIME itself can be further investigated by considering its different parameters, for example, the assignment of weights of perturbed examples in different ways, the robustness of the method to outliers, etc. The original Cox model has several modifications based on the Lasso method. These modifications could improve the explanation method in some applications, for example, in medicine applications to take into account a high dimensionality of survival data. This is also a direction for further research.

Acknowledgement

The reported study was funded by RFBR, project number 20-01-00154.

References

- [1] A. Adadi and M. Berrada. Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018.
- [2] I. Ahern, A. Noack, L. Guzman-Nateras, D. Dou, B. Li, and J. Huan. NormLime: A new feature importance metric for explaining deep neural networks. arXiv:1909.04200, Sep 2019.
- [3] A.B. Arrieta, N. Diaz-Rodriguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera. Explainable artificial in-

telligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. arXiv:1910.10045, October 2019.

- [4] V. Arya, R.K.E. Bellamy, P.-Y. Chen, A. Dhurandhar, M. Hind, S.C. Hoffman, S. Houde, Q.V. Liao, R. Luss, A. Mojsilovic, S. Mourad, P. Pedemonte, R. Raghavendra, J. Richards, P. Sattigeri, K. Shanmugam, M. Singh, K.R. Varshney, D. Wei, and Y. Zhang. One explanation does not fit all: A toolkit and taxonomy of AI explainability techniques. arXiv:1909.03012, Sep 2019.
- [5] F. Barthe, O. Guedon, S. Mendelson, and A. Naor. A probabilistic approach to the geometry of the l-ball. *The Annals of Probability*, 33(2):480–513, 2005.
- [6] V. Van Belle, K. Pelckmans, S. Van Huffel, and J.A. Suykens. Support vector methods for survival analysis: a comparison between ranking and regression approaches. *Artificial intelligence in medicine*, 53(2):107–118, 2011.
- [7] R. Bender, T. Augustin, and M. Blettner. Generating survival times to simulate cox proportional hazards models. *Statistics in medicine*, 24(11):1713–1723, 2005.
- [8] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [9] C. Burns, J. Thomason, and W. Tansey. Interpreting black box models with statistical guarantees. arXiv:1904.00045, Mar 2019.
- [10] D.V. Carvalho, E.M. Pereira, and J.S. Cardoso. Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(832):1–34, 2019.
- [11] D.R. Cox. Regression models and life-tables. *Journal of the Royal Statistical Society, Series B (Methodological)*, 34(2):187–220, 1972.
- [12] M. Du, N. Liu, and X. Hu. Techniques for interpretable machine learning. arXiv:1808.00033, May 2019.
- [13] D. Faraggi and R. Simon. A neural network model for survival data. *Statistics in medicine*, 14(1):73–82, 1995.
- [14] R. Fong and A. Vedaldi. Explanations for attributing deep neural network predictions. In *Explainable AI*, volume 11700 of *LNCS*, pages 149–167. Springer, Cham, 2019.
- [15] R.C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3429–3437. IEEE, 2017.
- [16] D. Garreau and U. von Luxburg. Explaining the explainer: A first theoretical analysis of LIME. arXiv:2001.03447, January 2020.
- [17] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee. Counterfactual visual explanations. arXiv:1904.07451, Apr 2019.
- [18] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi. A survey of methods for explaining black box models. *ACM computing surveys*, 51(5):93, 2019.
- [19] C. Haarburger, P. Weitz, O. Rippel, and D. Merhof. Image-based survival analysis for lung cancer patients using CNNs. arXiv:1808.09679v1, Aug 2018.

- [20] R. Harman and V. Lacko. On decompositional algorithms for uniform sampling from n-spheres and n-balls. *Journal of Multivariate Analysis*, 101:2297–2304, 2010.
- [21] F. Harrell, R. Califf, D. Pryor, K. Lee, and R. Rosati. Evaluating the yield of medical tests. *Journal of the American Medical Association*, 247:2543–2546, 1982.
- [22] L.A. Hendricks, R. Hu, T. Darrell, and Z. Akata. Grounding visual explanations. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 264–279, 2018.
- [23] D. Hosmer, S. Lemeshow, and S. May. *Applied Survival Analysis: Regression Modeling of Time to Event Data*. John Wiley & Sons, New Jersey, 2008.
- [24] L. Hu, J. Chen, V.N. Nair, and A. Sudjianto. Locally interpretable models and effects based on supervised partitioning (LIME-SUP). arXiv:1806.00663, Jun 2018.
- [25] Q. Huang, M. Yamada, Y. Tian, D. Singh, D. Yin, and Y. Chang. GraphLIME: Local interpretable model explanations for graph neural networks. arXiv:2001.06216, January 2020.
- [26] N.A. Ibrahim, A. Kudus, I. Daud, and M.R. Abu Bakar. Decision tree for competing risks survival probability in breast cancer study. *International Journal Of Biological and Medical Research*, 3(1):25–29, 2008.
- [27] S. Kaneko, A. Hirakawa, and C. Hamada. Enhancing the lasso approach for developing a survival prediction model based on gene expression data. *Computational and Mathematical Methods in Medicine*, 2015(Article ID 259474):1–7, 2015.
- [28] J.L. Katzman, U. Shaham, A. Cloninger, J. Bates, T. Jiang, and Y. Kluger. DeepSurv: Personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC medical research methodology*, 18(24):1–12, 2018.
- [29] F.M. Khan and V.B. Zubek. Support vector regression for censored data (SVRc): a novel tool for survival analysis. In *2008 Eighth IEEE International Conference on Data Mining*, pages 863–868. IEEE, 2008.
- [30] J. Kim, I. Sohn, S.-H. Jung, S. Kim, and C. Park. Analysis of survival data with group lasso. *Communications in Statistics - Simulation and Computation*, 41(9):1593–1605, 2012.
- [31] P.W. Koh and P. Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 1885–1894, 2017.
- [32] A. Van Looveren and J. Klaise. Interpretable counterfactual explanations guided by prototypes. arXiv:1907.02584, Jul 2019.
- [33] S.M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- [34] U.B. Mogensen, H. Ishwaran, and T.A. Gerds. Evaluating random forests for survival analysis using prediction error curves. *Journal of Statistical Software*, 50(11):1–23, 2012.
- [35] C. Molnar. *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*. Published online, <https://christophm.github.io/interpretable-ml-book/>, 2019.

- [36] W.J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yua. Interpretable machine learning: definitions, methods, and applications. *arXiv:1901.04592*, Jan 2019.
- [37] J.B. Nasejje, H. Mwambi, K. Dheda, and M. Lesosky. A comparison of the conditional inference survival forest model to random survival forests based on a simulation study as well as on two applications with time-to-event data. *BMC Medical Research Methodology*, 17(115):1–17, 2017.
- [38] I.K. Omurlu, M. Ture, and F. Tokatli. The comparisons of random survival forests and cox regression analysis with simulation and an application related to breast cancer. *Expert Systems with Applications*, 36:8582–8588, 2009.
- [39] V. Petsiuk, A. Das, and K. Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv:1806.07421*, June 2018.
- [40] J. Rabold, H. Deininger, M. Siebers, and U. Schmid. Enriching visual with verbal explanations for relational concepts – combining LIME with Aleph. *arXiv:1910.01837v1*, October 2019.
- [41] Y. Ramon, D. Martens, F. Provost, and T. Evgeniou. Counterfactual explanation algorithms for behavioral and textual data. *arXiv:1912.01819*, December 2019.
- [42] R. Ranganath, A. Perotte, N. Elhadad, and D. Blei. Deep survival analysis. *arXiv:1608.02158*, September 2016.
- [43] M.T. Ribeiro, S. Singh, and C. Guestrin. “Why should I trust you?” Explaining the predictions of any classifier. *arXiv:1602.04938v3*, Aug 2016.
- [44] M.T. Ribeiro, S. Singh, and C. Guestrin. Anchors: High-precision model-agnostic explanations. In *AAAI Conference on Artificial Intelligence*, pages 1527–1535, 2018.
- [45] C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1:206–215, 2019.
- [46] M. Schmid, M.N. Wright, and A. Ziegler. On the use of harrell’s c for clinical risk prediction via random survival forests. *Expert Systems with Applications*, 63:450–459, 2016.
- [47] S.M. Shankaranarayana and D. Runje. ALIME: Autoencoder based approach for local interpretability. *arXiv:1909.02437*, Sep 2019.
- [48] P.K. Shivaswamy, W. Chu, and M. Jansche. A support vector approach to censored targets. In *Seventh IEEE International Conference on Data Mining, ICDM 2007*, pages 655–660. IEEE, 2007.
- [49] E. Strumbel and I. Kononenko. An efficient explanation of individual classifications using game theory. *Journal of Machine Learning Research*, 11:1–18, 2010.
- [50] N. Ternes, F. Rotolo, and S. Michiels. Empirical extensions of the lasso penalty to reduce the false discovery rate in high-dimensional cox regression models. *Statistics in medicine*, 35(15):2561–2573, 2016.
- [51] R. Tibshirani. The lasso method for variable selection in the cox model. *Statistics in medicine*, 16(4):385–395, 1997.

- [52] J. van der Waa, M. Robeer, J. van Diggelen, M. Brinkhuis, and M. Neerincx. Contrastive explanations with local foil trees. arXiv:1806.07470, June 2018.
- [53] M.N. Vu, T.D. Nguyen, N. Phan, and M.T. Thai R. Gera. Evaluating explainers via perturbation. arXiv:1906.02032v1, Jun 2019.
- [54] S. Wachter, B. Mittelstadt, and C. Russell. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31:841–887, 2017.
- [55] H. Wang and L. Zhou. Random survival forest with space extensions for censored data. *Artificial intelligence in medicine*, 79:52–61, 2017.
- [56] P. Wang, Y. Li, and C.K. Reddy. Machine learning for survival analysis: A survey. arXiv:1708.04649, August 2017.
- [57] A. White and A.dA. Garcez. Measurable counterfactual local explanations for any classifier. arXiv:1908.03020v2, November 2019.
- [58] A. Widodo and B.-S. Yang. Machine health prognostics using survival probability and support vector machine. *Expert Systems with Applications*, 38(7):8430–8437, 2011.
- [59] D.M. Witten and R. Tibshirani. Survival analysis with high-dimensional covariates. *Statistical Methods in Medical Research*, 19(1):29–51, 2010.
- [60] M.N. Wright, T. Dankowski, and A. Ziegler. Random forests for survival analysis using maximally selected rank statistics. arXiv:1605.03391v1, May 2016.
- [61] M.N. Wright, T. Dankowski, and A. Ziegler. Unbiased split variable selection for random survival forests using maximally selected rank statistics. *Statistics in Medicine*, 36(8):1272–1284, 2017.
- [62] M.R. Zafar and N.M. Khan. DLIME: A deterministic local interpretable model-agnostic explanations approach for computer-aided diagnosis systems. arXiv:1906.10263, Jun 2019.
- [63] H.H. Zhang and W. Lu. Adaptive Lasso for Cox’s proportional hazards model. *Biometrika*, 94(3):691–703, 2007.
- [64] X. Zhu, J. Yao, and J. Huang. Deep convolutional neural network for survival analysis with pathological images. In *2016 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 544–547. IEEE, 2016.