

# Censor-aware Semi-supervised Learning for Survival Time Prediction from Medical Images <sup>\*</sup>

Renato Hermoza<sup>†</sup>    Gabriel Maicas<sup>†</sup>    Jacinto C. Nascimento<sup>‡</sup>  
Gustavo Carneiro<sup>†</sup>

<sup>†</sup>Australian Institute for Machine Learning, The University of Adelaide

<sup>‡</sup>Institute for Systems and Robotics, Instituto Superior Tecnico, Portugal

**Abstract.** Survival time prediction from medical images is important for treatment planning, where accurate estimations can improve health-care quality. One issue affecting the training of survival models is censored data. Most of the current survival prediction approaches are based on Cox models that can deal with censored data, but their application scope is limited because they output a hazard function instead of a survival time. On the other hand, methods that predict survival time usually ignore censored data, resulting in an under-utilization of the training set. In this work, we propose a new training method that predicts survival time using all censored and uncensored data. We propose to treat censored data as samples with a lower-bound time to death and estimate pseudo labels to semi-supervise a censor-aware survival time regressor. We evaluate our method on pathology and x-ray images from the TCGA-GM and NLST datasets. Our results establish the state-of-the-art survival prediction accuracy on both datasets.

**Keywords:** Censored data · Noisy labels · Pathological images · Chest x-rays · Semi-supervised Learning · Survival time prediction.

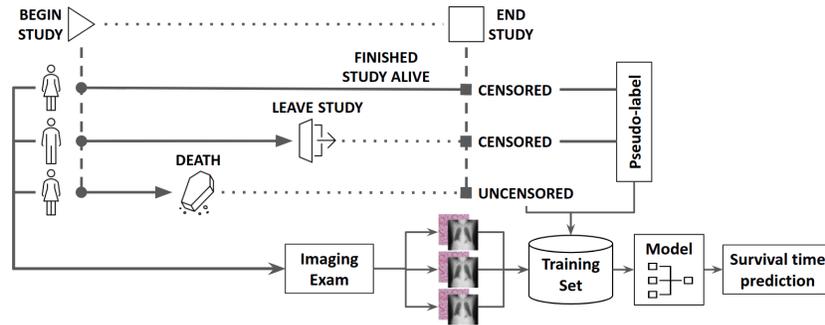
## 1 Introduction

Survival time prediction models estimate the time elapsed from the start of a study until an event (e.g., death) occurs with a patient. These models are important because they may influence treatment decisions that affect the health outcomes of patients [3]. Thus, automated models that produce accurate survival time estimations can be beneficial to improve the quality of healthcare.

Survival prediction analysis requires the handling of right-censored samples, which represent the cases where the event of interest has not occurred either because the study finished before the event happened, or the patient left the study before its end (Fig. 1). Current techniques to deal with censored data are typically based on Cox proportional hazards models [4] that rank patients in terms of their risks instead of predicting survival time, limiting their usefulness

---

<sup>\*</sup> This work was supported by the Australian Research Council through grants DP180103232 and FT190100525.



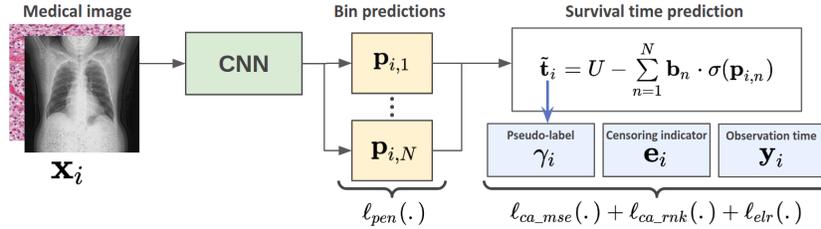
**Fig. 1.** At the beginning of the study, all patients are scanned for a pathology or chest x-ray image, and then patients are monitored until the end of study. Patients who die during the study represent uncensored data, while patients who do not die or leave before the end of study denote right-censored data. We train our model with uncensored and pseudo-labeled censored data to semi-supervise a censor-aware regressor to predict survival time.

for general clinical applications [2]. The problem is that the survival time prediction from a Cox hazard model requires the estimation of a baseline hazard function, which is challenging to obtain in practice [21]. A similar issue is observed in [17] that treats survival analysis as a ranking problem without directly estimating survival times.

Some methods [1,7,19] directly predict survival time, but they ignore the censored data for training. Other methods [8,21,20] use censored data during training, but in a sub-optimally manner. More specifically, these models use the censoring time as a lower bound for the event of interest. Even though this is better than disregarding the censored cases, such approaches do not consider the potential differences between the censoring time and the hidden survival time. We argue that if the values of these differences can be estimated with pseudo-labeling mechanisms, survival prediction models could be more accurate. However, pseudo labels estimated for the censored data may introduce noisy labels for training, which have been studied for classification [11,12] and segmentation [22], but never for survival prediction in a semi-supervision context.

In this paper, we propose a new training method for deep learning models to predict survival time from medical images using all censored and uncensored training data. The main contributions of this paper are:

- a method that estimates a pseudo label for the survival time of censored data (lower bounded by the annotated censoring time) that is used to semi-supervise a censor-aware survival time regressor; and
- two new regularization losses to handle pseudo label noise: a) we adapt the early-learning regularization (ELR) [12] loss from classification to survival prediction; and b) inspired by risk prediction models, we use a censor-aware ranking loss to produce a correct sorting of samples in terms of survival time.



**Fig. 2.** The proposed model outputs a set of  $N$  bin predictions  $\{\mathbf{p}_{i,n}\}_{n=1}^N$  (each bin representing an amount of time  $\mathbf{b}_n$ ) that are aggregated to produce a survival time prediction for the  $i^{th}$  case. This prediction is achieved by taking the maximum survival time  $U$  and subtracting it by the activation of each bin  $\sigma(p_{i,n})$  times the amount of time in  $\mathbf{b}_n$ . A set of loss functions are used while training: a censored-aware version of MSE ( $\ell_{ca\_mse}$ ), a penalization term for bin consistency ( $\ell_{pen}$ ), a rank loss ( $\ell_{ca\_rnk}$ ), and an adapted version of the ELR regularization [12] to survival prediction ( $\ell_{elr}$ ).

We evaluate our method using the TCGA-GM [13] dataset of pathological images, and the NLST dataset [14,15] of chest x-ray images. Our results show a clear benefit of using pseudo labels to train survival prediction models, achieving the new state-of-the-art (SOTA) survival time prediction results for both datasets. We make our code available at <https://github.com/renato145/CASurv>.

## 2 Method

To explain our method, we define the data set as  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{e}_i)\}_{i=1}^{|\mathcal{D}|}$ , where  $\mathbf{x}_i \in \mathcal{X}$  denotes a medical image with  $\mathcal{X} \subset \mathbb{R}^{H \times W \times C}$  ( $H, W, C$  denote image height, width and number of color channels),  $\mathbf{y}_i \in \mathbb{N}$  indicates the observation time defined in days, and  $\mathbf{e}_i \in \{0, 1\}$  is the censoring indicator. When  $\mathbf{e}_i = 0$  (uncensored observation),  $\mathbf{y}_i$  corresponds to a survival time  $\mathbf{t}_i$  which indicates the event of death for the individual. In cases where  $\mathbf{e}_i = 1$  (censored observation)  $\mathbf{t}_i$  is unknown, but is lower bounded by  $\mathbf{y}_i$  (i.e.  $\mathbf{t}_i > \mathbf{y}_i$ ).

### 2.1 Model

The architecture of our model extends the implementation from [7] to work with censored data (Fig. 2). Our survival time prediction model uses a Convolutional Neural Network (CNN) [10] represented by

$$f : \mathcal{X} \times \theta \rightarrow \mathcal{P}, \tag{1}$$

which is parameterized by  $\theta \in \Theta$  and takes an image  $\mathbf{x}$  to output a survival time confidence vector  $\mathbf{p} \in \mathcal{P}$  that represents the confidence on regressing to the number of days in each of the  $N$  bins of  $\mathcal{P} \subset \mathbb{R}^N$  that discretize  $\mathbf{y}_i$ . Each bin of the vector  $\mathbf{b} \in \mathbb{N}^N$  represents a number of days interval and the sum of them is the upper limit for the survival time in the dataset:  $U = \sum_{n=1}^N \mathbf{b}(n)$ . Bins

are discretized non-uniformly to balance number of samples per bin to avoid performance degradation [21].

To obtain the survival time prediction we first calculate the survival number of days per bin  $n$ , as in  $\tilde{\mathbf{p}}(n) = \mathbf{b}(n) \cdot \sigma(\mathbf{p}(n))$ , where  $\sigma(\cdot)$  represents the sigmoid function. The survival time prediction is obtained with  $\tilde{\mathbf{t}} = h(\mathbf{x}; \theta) = U - \sum_{n=1}^N \tilde{\mathbf{p}}(n)$ , where  $h(\mathbf{x}; \theta)$  represents the full survival time regression model. Hence, a larger activation  $\mathbf{p}(n)$  indicates higher risk, as the predicted  $\tilde{\mathbf{t}}$  will be smaller. The training of our model minimizes the loss function:

$$\begin{aligned} \ell(\mathcal{D}, \theta) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{e}_i, \mathbf{s}_i) \in \mathcal{D}} & [\ell_{ca\_mse}(\tilde{\mathbf{t}}_i, \mathbf{y}_i, \mathbf{e}_i, \mathbf{s}_i) + \alpha \cdot \ell_{elr}(\mathbf{p}_i) + \beta \cdot \ell_{pen}(\mathbf{p}_i) \\ & + \frac{1}{|\mathcal{D}|} \sum_{j=1}^{|\mathcal{D}|} \delta \cdot \ell_{ca\_rnk}(\tilde{\mathbf{t}}_i, \tilde{\mathbf{t}}_j, \mathbf{y}_i, \mathbf{y}_j, \mathbf{e}_i, \mathbf{e}_j)] \end{aligned} \quad (2)$$

where  $\ell_{ca\_mse}(\cdot)$  is a censor-aware mean squared error defined in (4),  $\mathbf{s}_i$  indicates if sample  $i$  is pseudo-labeled,  $\ell_{elr}(\cdot)$  is the ELR regularization [12] adapted for survival time prediction defined in (5),  $\ell_{pen}(\cdot)$  is the penalization term defined in (6),  $\ell_{ca\_rnk}(\cdot)$  is a censor-aware rank loss defined in (7), and  $\alpha$ ,  $\beta$  and  $\delta$  control the strength of  $\ell_{elr}(\cdot)$ ,  $\ell_{pen}(\cdot)$  and  $\ell_{ca\_rnk}(\cdot)$  losses.

**Pseudo Labels** We proposed the use of pseudo labels to semi-supervise our censor-aware survival time regressor. For a censored observation  $i$ , a pseudo label  $\gamma_i$  is estimated as:

$$\gamma_i = \max(\mathbf{y}_i, \tilde{\mathbf{t}}_i), \quad (3)$$

which takes advantage of the nature of censored data, as being lower bounded by  $\mathbf{y}_i$ . We estimate pseudo labels at the beginning of each epoch and treat them as uncensored observations by re-labeling  $\mathbf{y}_i = \gamma_i$  and setting  $\mathbf{s}_i = 1$ .

The quality of generated pseudo labels depends on the training procedure stage, where pseudo labels produced during the first epochs are less accurate than the ones at the last epochs. Hence, we control the ratio of censored sample labels to be replaced with pseudo labels at each epoch using a cosine annealing schedule to have no pseudo labels at the beginning of the training and all censored data with their pseudo labels by the end of the training.

**Censor-Aware Mean Squared Error (CA-MSE)** Using a regular MSE loss is not suitable for survival prediction because the survival time for censored observations is unknown, but lower-bounded by  $\mathbf{y}_i$ . Therefore, we assume an error exists for censored samples when  $\tilde{\mathbf{t}}_i < \mathbf{y}_i$ . However, when  $\tilde{\mathbf{t}}_i > \mathbf{y}_i$ , it will be incorrect to assume an error, as the real survival time is unknown. To mitigate this issue, we introduce the CA-MSE loss:

$$\ell_{ca\_mse}(\tilde{\mathbf{t}}_i, \mathbf{y}_i, \mathbf{e}_i, \mathbf{s}_i) = \begin{cases} 0, & \text{if } (\mathbf{e}_i = 1) \wedge (\tilde{\mathbf{t}}_i > \mathbf{y}_i) \\ \tau^{\mathbf{s}_i} (\tilde{\mathbf{t}}_i - \mathbf{y}_i)^2, & \text{otherwise} \end{cases}, \quad (4)$$

where  $\tau \in (0, 1)$  reduces the weight for pseudo-labeled data ( $\mathbf{s}_i = 1$ ).

**Early-Learning Regularization (ELR)** To deal with noisy pseudo labels, we modify ELR [12] to work in our survival time prediction setting, as follows:

$$\ell_{elr}(\mathbf{p}_i) = \log(1 - (1/N) (\sigma(\mathbf{p}_i)^\top \sigma(\mathbf{q}_i))), \quad (5)$$

where  $\mathbf{q}_i^{(k)} = \psi \mathbf{q}_i^{(k-1)} + (1 - \psi) \mathbf{p}_i^{(k)}$  (with  $\mathbf{q}_i \in \mathcal{P}$ ) is the temporal ensembling momentum [12] of the survival predictions, with  $k$  denoting the training epoch, and  $\psi \in [0, 1]$ . The idea of the loss in (5) is to never stop training for samples where the model prediction coincides with the temporal ensembling momentum (i.e., the ‘clean’ pseudo-labeled samples), and to never train for the noisy pseudo-labeled samples [12].

**Bin Penalization Term** We add the penalization term described in [7] that forces a bin  $n$  to be active only when all previous bins (1 to  $n - 1$ ) are active, forcing bins to represent sequential risk levels:

$$\ell_{pen}(\mathbf{p}_i) = \frac{1}{N-1} \sum_{n=1}^{N-1} \max(0, (\sigma(\mathbf{p}_i(n+1)) - \sigma(\mathbf{p}_i(n))))). \quad (6)$$

**Censor-Aware Rank Loss** Cox proportional hazards models [4] aim to rank training samples instead of the directly estimating survival time. Following this reasoning, we propose a censor-aware ranking loss to encourage this sample sorting behavior. The rank loss has the following form:

$$\ell_{ca\_rnk}(\tilde{\mathbf{t}}_i, \tilde{\mathbf{t}}_j, \mathbf{y}_i, \mathbf{y}_j, \mathbf{e}_i, \mathbf{e}_j) = \max(0, -\mathcal{G}(\mathbf{y}_i, \mathbf{y}_j, \mathbf{e}_i, \mathbf{e}_j) \times (\tilde{\mathbf{t}}_i - \tilde{\mathbf{t}}_j)) \quad (7)$$

where  $i$  and  $j$  index all pairs of samples in a training mini-batch, and

$$\mathcal{G}(\mathbf{y}_i, \mathbf{y}_j, \mathbf{e}_i, \mathbf{e}_j) = \begin{cases} +1, & \text{if } (\mathbf{y}_i > \mathbf{y}_j) \wedge \mathbf{e}_j = 0 \\ -1, & \text{if } (\mathbf{y}_i \leq \mathbf{y}_j) \wedge \mathbf{e}_i = 0 \\ 0, & \text{otherwise} \end{cases} . \quad (8)$$

In (8), when  $\mathcal{G}(\cdot) = +1$  the sample  $i$  should be ranked higher than sample  $j$ , when  $\mathcal{G}(\cdot) = -1$  the sample  $j$  should be ranked higher than sample  $i$ , and the loss is ignored when  $\mathcal{G}(\cdot) = 0$ .

### 3 Experiments

In this section, we describe the datasets and explain the experimental setup and evaluation measures, followed by the presentation of results.

#### 3.1 Data Sets

**TCGA-GM dataset** The Cancer Genome Atlas (TCGA) Lower-Grade Glioma (LGG) and Glioblastoma (GBM) projects: TCGA-GM [13] consists of 1,505 patches extracted from 1,061 whole-slide pathological images. A total of 769 unique patients with gliomas were included in the study with 381 censored and 388 uncensored observations. The labels are given in days and can last from 1 to 6423 days. We use the published train-test split provided with the data set [13].

**NLST dataset** The National Lung Screening Trial (NLST) [14,15] is a randomized multicenter study for early detection of lung cancer of current or former heavy smokers of ages 55 to 74. We only used the chest x-ray images from the dataset and excluded cases where the patient died for causes other than lung cancer. The original study includes 25,681 patients (77,040 images), and after filtering, we had 15,244 patients (47,947 images) with 272 uncensored cases. The labels are given in days and can last from 1 to 2983 days. We patient-wise split the dataset into training (9,133 patients), validation (3,058 patients) and testing (3,053 patients) sets, maintaining the same demographic distribution across sets.

### 3.2 Experimental Set Up

The input image size is adjusted to  $512^2$  pixels for both datasets and normalized with ImageNet [18] mean and standard deviation. For the model in (1), the space  $\mathcal{P}$  has 5 bins for TCGA-GM pathology images and 3 bins for NLST chest x-ray images. The model  $f(\mathbf{x}; \theta)$  in (1) is an ImageNet pre-trained ResNet-18 [6]. The training of the model uses Adam [9] optimizer with a momentum of 0.9, weight decay of 0.01 and a mini-batch size of 32 for 40 epochs. To obtain pseudo labels we have  $\tau = 0.5$  in (4), the bin penalization term  $\beta = 1e6$  in (6), the ELR [12] loss in (5) with  $\alpha = 100$  and  $\psi = 0.5$ , and the rank loss in (7) with  $\delta = 1$ . We run all experiments using PyTorch [16] v1.8 on a NVIDIA V100 GPU. The total training time is 45 minutes for the TCGA-GM dataset and 240 minutes for the NLST dataset, and the inference time for a single x-ray image takes 2.8 ms.

To evaluate our method, we use the Mean Absolute Error (MAE) between our predictions and the observed time, where the error is reported patient-wise, aggregating the mean of all image predictions per patient. To account for censored data, we use the following metric proposed by Xiao *et al.* [21].

$$\text{MAE} = \mathbb{E} [|\tilde{\mathbf{t}} - \mathbf{y}| \mathbb{I}(\tilde{\mathbf{t}} < \mathbf{y}) + (1 - \mathbf{e})|\tilde{\mathbf{t}} - \mathbf{y}| \mathbb{I}(\tilde{\mathbf{t}} \geq \mathbf{y})], \quad (9)$$

where  $\mathbb{I}$  is the indicator function. In (9), the censored cases where the model predicts the survival time  $\tilde{\mathbf{t}} \geq \mathbf{y}$  are not counted as errors. As a result, a higher count of censored records will decrease the average MAE, so we also measure the MAE considering only uncensored records. We also report concordance index (C-index) [5] to measure the correct ranking of the predictions.

### 3.3 Results

In Table 1, we compare survival time prediction results between our method and SOTA methods [13,21,23] on TCGA-GM, and DeepConvSurv [23] on NLST. Our method sets the new SOTA results in terms of MAE (all samples and only uncensored ones) and C-index for both datasets. The low MAE result for our approach on NLST can be explained by the fact that the vast majority of cases are censored (over 95%), and our method successfully estimates a survival time larger than the censoring time, producing MAE = 0 for many censored cases, according to (9).

**Table 1.** Comparison of survival time prediction between our approach and SOTA methods on the TCGA-GM and NLST. Evaluation metrics are mean absolute error (MAE) in days, and concordance index (C-index). Best results are highlighted.

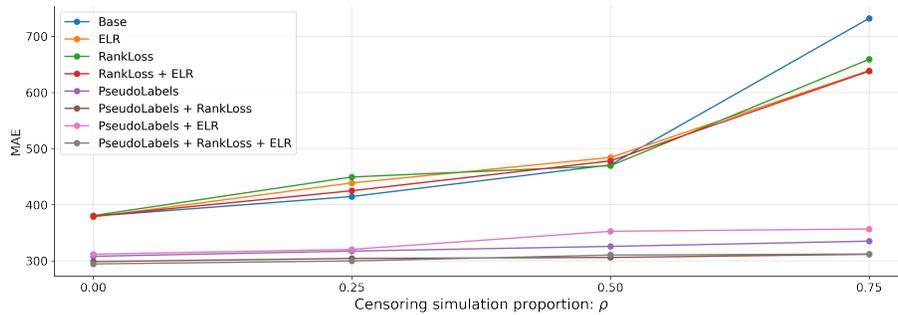
Method	MAE (all samples)	MAE (only uncensored)	C-index
TCGA-GM			
DeepConvSurv [23]	439.1	-	0.731
AFT [20]	386.5	-	0.685
SCNNs [13]	424.5	-	0.725
CDOR [21]	321.2	-	0.737
Ours	<b>286.95</b>	<b>365.06</b>	<b>0.740</b>
NLST			
DeepConvSurv [23]	27.98	1334.22	0.690
Ours	<b>26.28</b>	<b>1275.24</b>	<b>0.756</b>

**Table 2.** Ablation study for our method on the TCGA-GM and NLST datasets. The evaluation metrics are: mean absolute error (MAE) in days and concordance index (C-index). Best results are highlighted in bold.

$\ell_{ca\_mse} + \ell_{pen}$	$\ell_{ca\_rnk}$	Pseudo labels	$\ell_{elr}$	MAE (all samples)	Best MAE (all samples)	Best MAE (only uncensored)	C-index	Best C-index
TCGA-GM								
✓	-	-	-	379.40 ±20.67	356.74	521.91	0.702 ±0.009	0.707
✓	✓	-	-	380.49 ±17.55	361.76	517.41	0.722 ±0.008	0.730
✓	-	✓	-	308.01 ±7.99	302.04	405.48	0.725 ±0.014	<b>0.740</b>
✓	-	✓	✓	311.76 ±6.54	304.42	404.32	0.711 ±0.016	0.723
✓	✓	✓	-	298.56 ±6.31	291.53	411.79	0.718 ±0.016	0.736
✓	✓	✓	✓	<b>294.45 ±6.49</b>	<b>286.95</b>	<b>365.06</b>	<b>0.728 ±0.010</b>	<b>0.740</b>
NLST								
✓	-	-	-	27.24 ±0.15	27.10	1283.41	0.729 ±0.004	0.732
✓	✓	-	-	27.68 ±0.33	27.31	1286.27	0.732 ±0.003	0.736
✓	-	✓	-	<b>26.97 ±0.97</b>	<b>26.28</b>	<b>1275.24</b>	0.751 ±0.006	<b>0.756</b>
✓	-	✓	✓	29.05 ±0.59	28.37	1397.44	0.703 ±0.002	0.705
✓	✓	✓	-	27.83 ±2.32	26.29	1283.14	0.747 ±0.002	0.749
✓	✓	✓	✓	27.13 ±0.86	26.61	1277.02	<b>0.752 ±0.002</b>	0.754

On Table 2, we show the ablation study for the components explained in Section 2.1, where we start from a baseline model Base trained with  $\ell_{ca\_mse}$  and  $\ell_{pen}$ . Then we show how results change with the use of  $\ell_{ca\_rnk}$  (RankLoss), pseudo labels and  $\ell_{elr}$  (ELR). Regarding the MAE results on TCGA-GM and NLST datasets, we can observe a clear benefit of using pseudo labels. For TCGA-GM, the use of pseudo labels improve MAE by more than 20% compared with the Base model with RankLoss. For NLST, the smaller MAE reduction with pseudo labels can be explained by the large proportion of censored cases, which makes the robust estimation of survival time a challenging task. Considering the C-index, pseudo labels show similar benefits for TCGA-GM, with a slightly better result. On NLST, pseudo labels improve the C-index from 0.73 to 0.75 over the Base model with RankLoss. On both datasets, the use of Base model with RankLoss, pseudo labels, and ELR, presents either the best results or a competitive result with best result.

The differences between the results on NLST and TCGA-GM can be attributed to different proportions of censored/uncensored cases in both datasets.



**Fig. 3.** Performance of our method under different proportions of censored vs uncensored labels on the TCGA-GM dataset. The x-axis shows the ratio of uncensored records transformed to be censored ( $\rho$ ), and the y-axis shows MAE. Note that pseudo labels combined with RankLoss and ELR show robust results in scenarios with small percentages of uncensored records.

We test the hypothesis that such differences have a strong impact in the performance of our algorithm by formulating an experiment using the TCGA-GM dataset, where we simulate different proportion of censored records. Figure 3 shows the result of this simulation, where we show MAE as a function of the proportion of censored data, denoted by  $\rho$ . The results show that the PseudoLabels methods (combined with RankLoss and ELR losses) keep the MAE roughly constant, while Base, ELR and RankLoss results deteriorate significantly. The best results are achieved by PseudoLabels + RankLoss and PseudoLabels + RankLoss + ELR, suggesting that pseudo labels and both losses are important for the method to be robust to large proportions of censored data.

## 4 Discussion and Conclusion

Current techniques for survival prediction discard or sub-optimally use the variability present in censored data. In this paper we proposed a method to exploit such variability on censored data by using pseudo labels to semi-supervise the learning process. Experimental results in Table 1 showed that our proposed method obtains SOTA results for survival time prediction on the TCGA-GM and NLST datasets. In fact, the effect of pseudo labels can be observed in Figure 3, where we artificially censored samples from the TCGA-GM dataset. It is clear that the use of pseudo labels provide solid robustness to a varying proportion of censored data in the training set. Combining the pseudo labels with a regularisation loss that accounts for noisy pseudo labels and censor-aware rank loss improves this robustness and results on datasets. However, is important to note that the method may not be able to produce good quality pseudo-labels on cases where the dataset contains few uncensored records. We expect our newly proposed method to foster the development of new survival time prediction models that exploit the variability present in censored data.

## References

1. Agravat, R.R., Raval, M.S.: Brain Tumor Segmentation and Survival Prediction. In: International MICCAI Brainlesion Workshop. pp. 338–348. Springer (2019)
2. Baid, U., Rane, S.U., Talbar, S., Gupta, S., Thakur, M.H., Moiyadi, A., Mahajan, A.: Overall Survival Prediction in Glioblastoma With Radiomic Features Using Machine Learning. *Frontiers in Computational Neuroscience* **14** (2020). <https://doi.org/10.3389/fncom.2020.00061>, publisher: Frontiers
3. Cheon, S., Agarwal, A., Popovic, M., Milakovic, M., Lam, M., Fu, W., DiGiovanni, J., Lam, H., Lechner, B., Pulenzas, N., Chow, R., Chow, E.: The accuracy of clinicians' predictions of survival in advanced cancer: a review. *Annals of Palliative Medicine* **5**(1), 22–29 (Jan 2016). <https://doi.org/10.3978/j.issn.2224-5820.2015.08.04>
4. Cox, D.R.: Regression Models and Life-Tables. *Journal of the Royal Statistical Society. Series B (Methodological)* **34**(2), 187–220 (1972)
5. Harrell, F.E., Califf, R.M., Pryor, D.B., Lee, K.L., Rosati, R.A.: Evaluating the yield of medical tests. *Jama* **247**(18), 2543–2546 (1982), publisher: American Medical Association
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778 (Jun 2016). <https://doi.org/10.1109/CVPR.2016.90>
7. Hermoza, R., Maicas, G., Nascimento, J.C., Carneiro, G.: Post-hoc overall survival time prediction from brain mri. In: 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI). pp. 1476–1480 (2021). <https://doi.org/10.1109/ISBI48211.2021.9433877>
8. Jing, B., Zhang, T., Wang, Z., Jin, Y., Liu, K., Qiu, W., Ke, L., Sun, Y., He, C., Hou, D., Tang, L., lv, x., Li, C.: A deep survival analysis method based on ranking. *Artificial intelligence in medicine* **98**, 1–9 (07 2019). <https://doi.org/10.1016/j.artmed.2019.06.001>
9. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations (ICLR)* (2015)
10. Lecun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (Nov 1998). <https://doi.org/10.1109/5.726791>, conference Name: Proceedings of the IEEE
11. Li, J., Socher, R., Hoi, S.C.: Dividemix: Learning with noisy labels as semi-supervised learning. *arXiv preprint arXiv:2002.07394* (2020)
12. Liu, S., Niles-Weed, J., Razavian, N., Fernandez-Granda, C.: Early-Learning Regularization Prevents Memorization of Noisy Labels. *arXiv:2007.00151 [cs, stat]* (Oct 2020), <http://arxiv.org/abs/2007.00151>, arXiv: 2007.00151
13. Mobadersany, P., Yousefi, S., Amgad, M., Gutman, D.A., Barnholtz-Sloan, J.S., Vega, J.E.V., Brat, D.J., Cooper, L.A.D.: Predicting cancer outcomes from histology and genomics using convolutional networks. *Proceedings of the National Academy of Sciences* **115**(13), E2970–E2979 (Mar 2018). <https://doi.org/10.1073/pnas.1717139115>, <https://www.pnas.org/content/115/13/E2970>, publisher: National Academy of Sciences Section: PNAS Plus
14. National Lung Screening Trial Research Team, Aberle, D.R., Adams, A.M., Berg, C.D., Black, W.C., Clapp, J.D., Fagerstrom, R.M., Gareen, I.F., Gatsonis, C., Marcus, P.M., Sicks, J.D.: Reduced lung-cancer mortality with low-dose computed tomographic screening. *The New England journal of medicine* **365**(5), 395–409 (Aug 2011). <https://doi.org/10.1056/NEJMoa1102873>, <https://pubmed.ncbi.nlm.nih.gov/21714641>, edition: 2011/06/29

15. National Lung Screening Trial Research Team, Aberle, D.R., Berg, C.D., Black, W.C., Church, T.R., Fagerstrom, R.M., Galen, B., Gareen, I.F., Gatsonis, C., Goldin, J., Gohagan, J.K., Hillman, B., Jaffe, C., Kramer, B.S., Lynch, D., Marcus, P.M., Schnall, M., Sullivan, D.C., Sullivan, D., Zylak, C.J.: The National Lung Screening Trial: overview and study design. *Radiology* **258**(1), 243–253 (Jan 2011). <https://doi.org/10.1148/radiol.10091808>, <https://pubmed.ncbi.nlm.nih.gov/21045183>, edition: 2010/11/02 Publisher: Radiological Society of North America, Inc.
16. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: PyTorch: An Imperative Style, High-Performance Deep Learning Library. In: Wallach, H., Larochelle, H., Beygelzimer, A., Alché-Buc, F.d., Fox, E., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 32, pp. 8024–8035. Curran Associates, Inc. (2019)
17. Raykar, V.C., Steck, H., Krishnapuram, B., Dehing-Oberije, C., Lambin, P.: On ranking in survival analysis: bounds on the concordance index. In: *Proceedings of the 20th International Conference on Neural Information Processing Systems*. pp. 1209–1216. NIPS’07, Curran Associates Inc., Red Hook, NY, USA (Dec 2007)
18. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* **115**(3), 211–252 (Dec 2015). <https://doi.org/10.1007/s11263-015-0816-y>
19. Tang, Z., Xu, Y., Jiao, Z., Lu, J., Jin, L., Aibaidula, A., Wu, J., Wang, Q., Zhang, H., Shen, D.: Pre-operative Overall Survival Time Prediction for Glioblastoma Patients Using Deep Learning on Both Imaging Phenotype and Genotype. In: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (eds.) *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*. pp. 415–422. *Lecture Notes in Computer Science*, Springer International Publishing, Cham (2019). [https://doi.org/10.1007/978-3-030-32239-7\\_46](https://doi.org/10.1007/978-3-030-32239-7_46)
20. Wei, L.J.: The accelerated failure time model: a useful alternative to the cox regression model in survival analysis. *Statistics in medicine* **11**(14-15), 1871–1879 (1992)
21. Xiao, L., Yu, J.G., Liu, Z., Ou, J., Deng, S., Yang, Z., Li, Y.: Censoring-Aware Deep Ordinal Regression for Survival Prediction from Pathological Images. In: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (eds.) *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*. pp. 449–458. *Lecture Notes in Computer Science*, Springer International Publishing, Cham (2020). [https://doi.org/10.1007/978-3-030-59722-1\\_43](https://doi.org/10.1007/978-3-030-59722-1_43)
22. Zheng, Z., Yang, Y.: Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *International Journal of Computer Vision* pp. 1–15 (2021)
23. Zhu, X., Yao, J., Huang, J.: Deep convolutional neural network for survival analysis with pathological images. In: *2016 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. pp. 544–547 (Dec 2016). <https://doi.org/10.1109/BIBM.2016.7822579>

## A Supplementary Material

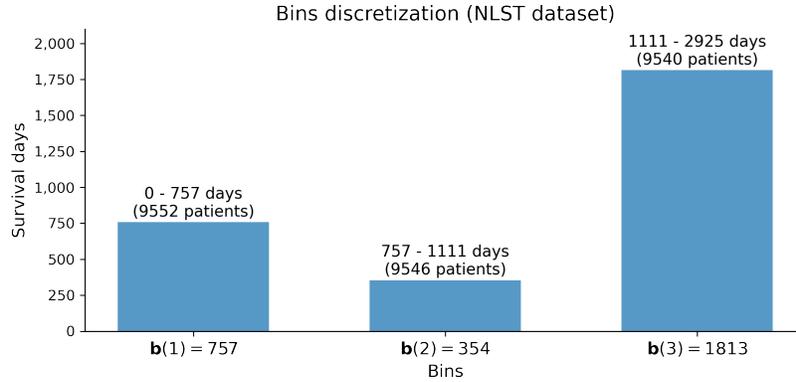
---

**Algorithm 1** Pseudocode for the training procedure

---

**Require:**  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{e}_i)\}_{i=1}^{|\mathcal{D}|}$  = training data  
**Require:**  $h(\mathbf{x}; \theta)$  = survival model with parameters  $\theta$   
**Require:**  $\alpha$  = ELR loss strength  
**Require:**  $\beta$  = Bin penalization strength  
**Require:**  $\delta$  = Rank loss strength  
**for**  $k$  in  $[1, k_{total}]$  **do**  
    *// Ratio of pseudo-labels from censored samples*  
     $m \leftarrow (1 + \cos \pi(1 - k/k_{total}))/2$   
    *// Set of censored samples to use pseudo-labels*  
     $P \leftarrow \text{random\_sample}(\mathcal{D}, m)$   
    **for**  $i$  in  $P$  **do**  
        *// Re-label censored sample*  
         $\tilde{\mathbf{t}}_i \leftarrow h(\mathbf{x}_i; \theta)$   
         $\mathbf{y}_i \leftarrow \max(\mathbf{y}_i, \tilde{\mathbf{t}}_i)$   
         $\mathbf{e}_i \leftarrow 1$   
    **end for**  
    *// Update network parameters*  
    **for** each minibatch  $B$  **do**  
        **for**  $j$  in  $B$  **do**  
             $\tilde{\mathbf{t}}_j \leftarrow h(\mathbf{x}_j; \theta)$   
        **end for**  
         $loss \leftarrow \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \ell_{ca\_mse}(\tilde{\mathbf{t}}_i, \mathbf{y}_i, \mathbf{e}_i)$   
             $+ \alpha \cdot \ell_{elr}(\mathbf{p}_i) + \beta \cdot \ell_{pen}(\mathbf{p}_i)$   
             $+ \frac{1}{|\mathcal{D}|} \sum_{j=1}^{|\mathcal{D}|} \delta \cdot \ell_{ca\_rnk}(\tilde{\mathbf{t}}_i, \tilde{\mathbf{t}}_j, \mathbf{y}_i, \mathbf{y}_j, \mathbf{e}_i, \mathbf{e}_j)$   
        update  $\theta$  using SGD with  $loss$   
    **end for**  
**end for**  
**return**  $\theta$

---



**Fig. S1.** We discretize the NLST dataset survival time into three bins, each representing a different amount of time, as depicted in the plot (e.g.,  $\mathbf{b}(1) = 757$  days,  $\mathbf{b}(2) = 354$  days, and  $\mathbf{b}(3) = 1813$  days). The discretization strategy aims to balance the number of patients per bin. Note that this discretization strategy uses only the information in the training set.

**Table S1.** One-sided Wilcoxon signed-rank test for the TCGA-GM dataset. We compare the mean absolute error (MAE) of our best performing method (Pseudolabels +  $\ell_{ca\_rnk} + \ell_{elr}$ ) against the rest. \* indicates that the best performing method is better assuming a significance level of 5%.

$\ell_{pen}$	$\ell_{ca\_rnk}$	Pseudo-labels	$\ell_{elr}$	Wilcox	p-value
✓	-	-	-	2157.0	0.0115 *
✓	✓	-	-	2069.0	0.0113 *
✓	-	✓	-	2383.0	0.0579
✓	-	✓	✓	2633.0	0.3164
✓	✓	✓	-	2680.0	0.3716

**Table S2.** One-sided Wilcoxon signed-rank test for the NLST dataset. We compare the mean absolute error (MAE) of our best performing method (PseudoLabels) against the rest. \* indicates that the best performing method is better assuming a significance level of 5%.

$\ell_{pen}$	$\ell_{ca\_rnk}$	Pseudo-labels	$\ell_{elr}$	Wilcox	p-value
✓	-	-	-	4651.0	0.0023 *
✓	✓	-	-	5252.0	0.0284 *
✓	-	✓	✓	2361.0	0.0099 *
✓	✓	✓	-	4807.0	0.2051
✓	✓	✓	✓	4061.0	0.0593

# Censor-aware Semi-supervised Learning for Survival Time Prediction from Medical Images

Paper 894

Anonymous

---

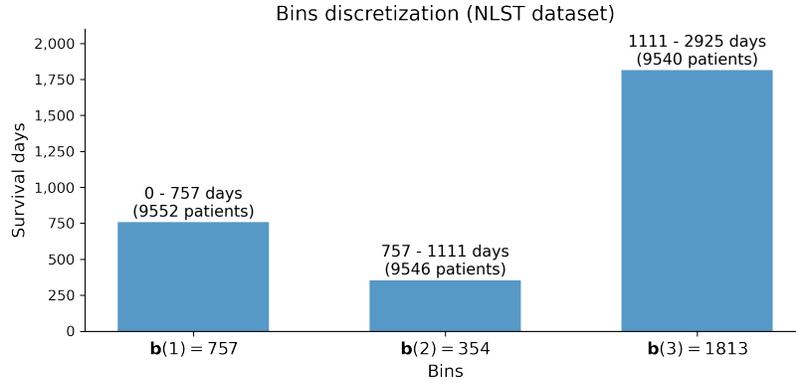
**Algorithm 1** Pseudocode for the training procedure

---

**Require:**  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{e}_i)\}_{i=1}^{|\mathcal{D}|}$  = training data  
**Require:**  $h(\mathbf{x}; \theta)$  = survival model with parameters  $\theta$   
**Require:**  $\alpha$  = ELR loss strength  
**Require:**  $\beta$  = Bin penalization strength  
**Require:**  $\delta$  = Rank loss strength

**for**  $k$  in  $[1, k_{total}]$  **do**  
  // Ratio of pseudo-labels from censored samples  
   $m \leftarrow (1 + \cos \pi(1 - k/k_{total}))/2$   
  // Set of censored samples to use pseudo-labels  
   $P \leftarrow \text{random\_sample}(\mathcal{D}, m)$   
  **for**  $i$  in  $P$  **do**  
    // Re-label censored sample  
     $\tilde{\mathbf{t}}_i \leftarrow h(\mathbf{x}_i; \theta)$   
     $\mathbf{y}_i \leftarrow \max(\mathbf{y}_i, \tilde{\mathbf{t}}_i)$   
     $\mathbf{e}_i \leftarrow 1$   
  **end for**  
  // Update network parameters  
  **for** each minibatch  $B$  **do**  
    **for**  $j$  in  $B$  **do**  
       $\tilde{\mathbf{t}}_j \leftarrow h(\mathbf{x}_j; \theta)$   
    **end for**  
     $loss \leftarrow \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \ell_{ca\_mse}(\tilde{\mathbf{t}}_i, \mathbf{y}_i, \mathbf{e}_i)$   
       $+ \alpha \cdot \ell_{elr}(\mathbf{p}_i) + \beta \cdot \ell_{pen}(\mathbf{p}_i)$   
       $+ \frac{1}{|\mathcal{D}|} \sum_{j=1}^{|\mathcal{D}|} \delta \cdot \ell_{ca\_rnk}(\tilde{\mathbf{t}}_i, \tilde{\mathbf{t}}_j, \mathbf{y}_i, \mathbf{y}_j, \mathbf{e}_i, \mathbf{e}_j)$   
    update  $\theta$  using SGD with  $loss$   
  **end for**  
**end for**  
**return**  $\theta$

---



**Fig. S1.** We discretize the NLST dataset survival time into three bins, each representing a different amount of time, as depicted in the plot (e.g.,  $\mathbf{b}(1) = 757$  days,  $\mathbf{b}(2) = 354$  days, and  $\mathbf{b}(3) = 1813$  days). The discretization strategy aims to balance the number of patients per bin. Note that this discretization strategy uses only the information in the training set.

**Table S1.** One-sided Wilcoxon signed-rank test for the TCGA-GM dataset. We compare the mean absolute error (MAE) of our best performing method (Pseudolabels +  $l_{ca\_rnk} + l_{etr}$ ) against the rest. \* indicates that the best performing method is better assuming a significance level of 5%.

$l_{pen}$	$l_{ca\_rnk}$	Pseudo-labels	$l_{etr}$	Wilcox	p-value
✓	-	-	-	2157.0	0.0115 *
✓	✓	-	-	2069.0	0.0113 *
✓	-	✓	-	2383.0	0.0579
✓	-	✓	✓	2633.0	0.3164
✓	✓	✓	-	2680.0	0.3716

**Table S2.** One-sided Wilcoxon signed-rank test for the NLST dataset. We compare the mean absolute error (MAE) of our best performing method (PseudoLabels) against the rest. \* indicates that the best performing method is better assuming a significance level of 5%.

$l_{pen}$	$l_{ca\_rnk}$	Pseudo-labels	$l_{etr}$	Wilcox	p-value
✓	-	-	-	4651.0	0.0023 *
✓	✓	-	-	5252.0	0.0284 *
✓	-	✓	✓	2361.0	0.0099 *
✓	✓	✓	-	4807.0	0.2051
✓	✓	✓	✓	4061.0	0.0593