

1 We thank the reviewers for their feedback. We will release our implementation on github. We thank R1 for pointing out
 2 our calibration metric of choice, D-Calibration (D-Cal), has been published at JMLR 2020 (we cite this as [H] here).

3 **[R3: Relation to D-Cal]** Building on [H]’s D-cal, we propose X-cal to regularize a model to have low D-Cal. X-cal is
 4 a differentiable approximation of an upper-bound on D-Cal, amenable to stochastic optimization.

5 **[R1, R2: Real world dataset with censoring; Survival benchmarks]** We evaluated X-Cal on [Avati et al.]’s alternative
 6 MIMIC set with 70% censored points. D-Cal goes from $2 \times 10^{-4} \rightarrow 9 \times 10^{-5}$ as λ increases from $0 \rightarrow 10^3$. We will
 7 add this to the paper. We are glad to include specific evaluations/benchmark sets that the reviewers think are relevant.

8 **[R2: Comparison with MTLR/ approaches mentioned in [H]]** We do not include methods from [H] because our
 9 work focuses on using flexible models with good likelihood but poor calibration, like S-CRPS. We use a categorical
 10 model because it is a common flexible likelihood that can approximate many continuous distributions given enough bins.
 11 This allows us to evaluate X-Cal without parametric restrictions. We did cite an approach like MTLR [Ranganath 2018].
 12 We will cite MTLR [Yu et al. 2011] and its neural version [Fotso 2018]. We ran MTLR using PyCox library on the
 13 uncensored synthetic gamma dataset. This gave a model with D-cal 0.7486, which is higher than any model we study.

14 **[R3,R4: Comparison to other established calibration metrics]** The alternative notion of calibration for fixed time t^*
 15 suggested by reviewers [Yadlowsky 2019, Royston/Altman 2013] are described in [H] as “1-Calibration”. [H] proves
 16 that D-Cal (with fixed bins) and 1-Cal for time t^* (with fixed bins) measure different aspects of the survival distribution:
 17 0 D-Cal and 0 1-Cal do not imply each other. A practitioner may need calibration at several times e.g. 6 months, 1 year.
 18 Future work is to regularize models with approximations of 1-Cal. measures (e.g. Hosmer-Lemeshow statistic) using
 19 soft indicators. Our focus is to maintain a certain level of calibration based on the specific metric, D-cal.

20 **[R2: p -value]** The p -value reported by [H] is the result of a χ^2 -test on the D-Cal test statistic. Thus, if models are
 21 ordered in the test statistic their corresponding p -values are ordered in the same way. While p -values help test for perfect
 22 calibration, our focus is on *improving* calibration of existing models which we demonstrate in in our experiments.

23 **[R3: Discontinuous learned conditional survival model]** As mentioned in 4.1, a discontinuous model will have a
 24 lower bound greater than 0 for D-Cal because its CDF values will not be a continuous Unif(0,1) variable. However,
 25 minimizing D-Cal will still spread out the CDF values to whichever extent possible and thus improve calibration. In the
 26 case of a categorical model, this lower bound decreases to 0 as the size of each bin goes to 0 when adding more bins.

27 **[R3: Adjustments for right censoring / MNIST censoring]** This is an important issue. In line 151 of our paper, we
 28 handle right censoring using the technique proposed in appendix B.5 in [H] and proved to result in a valid test statistic.
 29 As noted in [H] on page 47, the estimation of D-cal on a censored dataset will not equal the estimate when censored
 30 times are revealed. This is due to the fact that in the censored dataset [H]’s correction for right censoring gives a
 31 few bins the correct weight for free meaning D-cal will be lesser. However, for a given dataset, D-Cal is 0 for any
 32 bin for the true conditional $p(T | X)$ for any non-informative censoring process that meets a "positivity" assumption.
 33 Thus, two models evaluated on the *same* data (censored or uncensored) can be compared with D-Cal. Reweighting
 34 methods, such as Yadlowsky et al. that R3 suggests, can be used to adjust for censoring. One option is to adjust with
 35 with $p(C | X)$. This requires $C \perp T | X$ and solving a censored survival problem $p(C | X)$ with a high-dimensional
 36 conditioning set. Another option is to adjust with the lower dimensional conditioning set $p(C | risk_{\theta}(X))$. This
 37 requires $C \perp T | risk_{\theta}(X)$ and differentiating through the *estimation* of $p(C | risk_{\theta}(X))$ w.r.t. θ . The approach we
 38 take requires neither another censored survival problem nor differentiating through estimation.

39 **[R2: λ and γ]** Choosing λ : the user first decides on a
 40 threshold of D-Cal and then increase λ until D-Cal eval-
 41 uated on a held-out validation set meets this threshold.
 42 See Table 1 for the role of γ . With low γ , soft D-cal approximates
 43 poorly and D-Cal/NLL suffer. For γ too large, gradients vanish
 44 and D-Cal/NLL suffer. We found $\gamma = 10^4$ allowed for easy
 45 optimization with soft D-Cal approximating D-Cal well.

γ	10	10^4	1.1×10^4	10^5
D-Cal	0.4	0.0005	0.0002	0.0003
NLL	4.33	1.82	1.88	2.49

Table 1: Soft D-cal as γ varies.

46 **[R2, Tightness of upper-bound]** Table 2 shows that models ordered by
 47 the upper-bound are ordered in D-cal the same way. Further, when batch
 48 size is large enough, if $\lambda_i < \lambda_j$, the bound for λ_j is less than D-Cal for λ_i .

49 **[R4: Choosing D-Cal on a validation set]** During training we evaluated
 50 NLL+D-cal on a validation set at every epoch and save the model. Then,
 51 we report test metrics for the model with best validation NLL + D-Cal. If
 52 we only select/optimize X-cal on a validation set, the predictive likelihood
 53 may get arbitrarily worse. This issue occurs with Platt scaling as well.

54 **[Minor comments/ Typos]** We thank the reviewers for detailed feedback
 55 about our writing. We will define Harrell’s Concordance Index, change 10k
 56 to our intended 10^4 , and rephrase "calibration means accurate prognosis".

λ	Batch size	D-Cal	Bound
10	500	0.0040	0.0059
	5000	"	0.0042
50	500	0.0006	0.0024
	5000	"	0.0008
100	500	0.0003	0.0022
	5000	"	0.0005

Table 2: Slack in the Upper Bound