

# Neural Frailty Machine: Beyond proportional hazard assumption in neural survival regressions

Ruofan Wu\*, Jiawei Qiao\*, Mingzhe Wu, Wen Yu, Ming Zheng,  
Tengfei Liu, Tianyi Zhang, Weiqiang Wang

Ant Group, Fudan University and Coupang

- We study regression analysis with time-to-event data under **right censoring**.
- Event time  $\tilde{T} \geq 0$  with survival function  $S(t)$ , density  $f(t)$  and hazard function  $\lambda(t)$ .
- Feature vector  $Z$ .
- Censoring time  $C$  which satisfies  $\tilde{T} \perp\!\!\!\perp C|Z$ , resulting in observed tuple  $T = \tilde{T} \wedge C, \delta = I(\tilde{T} \leq C)$ .

## CoxPH model [Cox72]

In its linear form, the widely used CoxPH model assumes the proportional hazard (PH) assumption

$$\lambda(t|Z) = \lambda_0(t)e^{(Z, \theta)} \quad (1)$$

- Due to limited expressivity, there have been attempts that generalizes CoxPH using more elegant function approximators like neural networks [FS95, KSC<sup>+</sup>18].
- Nonlinear CoxPH models still assumes proportional hazard assumption which has limitations in modeling specific phenomenons like crossing hazards [Ben83].
- A more general solution: Using the idea of **frailty**.

## Frailty model [Hou86]

In its linear form, (multiplicative) frailty model introduces an additional unobserved heterogeneity  $\omega$  into the hazard/intensity formulation:

$$\lambda(t|Z, \omega) = \omega \lambda_0(t) e^{\langle Z, \theta \rangle} \quad (2)$$

To further enhance the capability of frailty model, we equip neural function approximations. The resulting modeling framework is called Neural Frailty Machine (NFM) which we propose to approaches.

## Proportional frailty scheme (PF)

The PF scheme directly replaces linear term in ordinary frailty model with a neural network  $m(Z)$  that depends only on the features.

$$\lambda(t|Z, \omega) = \omega e^{h(t)+m(Z)} \quad (3)$$

## Fully neural scheme (FN)

The FN scheme further relaxes the separation between baseline hazard and feature dependence, using a neural network  $h(t, Z)$  for approximation.

$$\lambda(t|Z, \omega) = \omega e^{\nu(t, Z)} \quad (4)$$

We use observed log-likelihood as the learning objective.

**PF-Scheme** We use two MLPs  $\hat{h} = \hat{h}(t; \mathbf{W}^h, \mathbf{b}^h)$  and  $\hat{m} = \hat{m}(Z; \mathbf{W}^m, \mathbf{b}^m)$  as function approximators to  $h$  and  $m$ .

$$\begin{aligned} & \mathcal{L}(\mathbf{W}^h, \mathbf{b}^h, \mathbf{W}^m, \mathbf{b}^m, \theta) \\ &= \frac{1}{n} \left[ \sum_{i \in [n]} \delta_i \log g_\theta \left( e^{\hat{m}(Z_i)} \int_0^{T_i} e^{\hat{h}(s)} ds \right) + \delta_i \hat{h}(T_i) + \delta_i \hat{m}(Z_i) - G_\theta \left( e^{\hat{m}(Z_i)} \int_0^{T_i} e^{\hat{h}(s)} ds \right) \right]. \end{aligned}$$

**FN-Scheme** We use  $\hat{\nu} = \hat{\nu}(t, Z; \mathbf{W}^\nu, \mathbf{b}^\nu)$  to approximate  $\nu(t, Z)$

$$\begin{aligned} & \mathcal{L}(\mathbf{W}^\nu, \mathbf{b}^\nu, \theta) \\ &= \frac{1}{n} \left[ \sum_{i \in [n]} \delta_i \log g_\theta \left( \int_0^{T_i} e^{\hat{\nu}(s, Z_i; \mathbf{W}^\nu, \mathbf{b}^\nu)} ds \right) + \delta_i \hat{\nu}(T_i, Z_i; \mathbf{W}^\nu, \mathbf{b}^\nu) - G_\theta \left( \int_0^{T_i} e^{\hat{\nu}(s, Z_i; \mathbf{W}^\nu, \mathbf{b}^\nu)} ds \right) \right]. \end{aligned}$$

Here  $G_\theta$  is defined as the negative of the logarithm of the Laplace transform of the frailty distribution, with  $g_\theta$  being its derivative w.r.t.  $t$ .

We study the rates of convergence for estimated parameters with the true parameters lying inside a Hölder ball with radius  $M$  and smoothness parameter  $\beta$ . Under the following two distance metrics

$$d_{PF}(\hat{\phi}_n, \phi_0) = \sqrt{\mathbb{E}_{Z \sim \mathbb{P}_Z} \left[ H^2(\mathbb{P}_{\hat{\phi}_n, Z=Z} \parallel \mathbb{P}_{\phi_0, Z=Z}) \right]}, \quad d_{FN}(\hat{\psi}_n, \psi_0) = \sqrt{\mathbb{E}_{Z \sim \mathbb{P}_Z} \left[ H^2(\mathbb{P}_{\hat{\psi}_n, Z=Z} \parallel \mathbb{P}_{\psi_0, Z=Z}) \right]}$$

where  $\hat{\phi}_n$  and  $\hat{\psi}_n$  are bundled parameter updates and  $H$  denotes Hellinger distance, we have the following statistical guarantee:

### Theorem

*Under some regularity conditions, we have*

$$d_{PF}(\hat{\phi}_n, \phi_0) = \tilde{O}_{\mathbb{P}} \left( n^{-\frac{\beta}{2\beta+2d}} \right), \quad d_{FN}(\hat{\psi}_n, \psi_0) = \tilde{O}_{\mathbb{P}} \left( n^{-\frac{\beta}{2\beta+2d+2}} \right)$$

*where  $d$  is the feature dimension and logarithmic factors are hidden.*

Empirical evaluations over 4 relatively small scale datasets.

Model	METABRIC		RotGBSG		FLCHAIN		SUPPORT	
	IBS	INBLL	IBS	INBLL	IBS	INBLL	IBS	INBLL
CoxPH	16.46 $\pm$ 0.90	49.57 $\pm$ 2.66	18.25 $\pm$ 0.44	53.76 $\pm$ 1.11	10.05 $\pm$ 0.38	33.18 $\pm$ 1.16	20.54 $\pm$ 0.38	59.58 $\pm$ 0.86
GBM	16.61 $\pm$ 0.82	49.87 $\pm$ 2.44	17.83 $\pm$ 0.44	52.78 $\pm$ 1.11	<u>9.98</u> $\pm$ 0.37	<u>32.88</u> $\pm$ 1.05	19.18 $\pm$ 0.39	56.46 $\pm$ 0.10
RSF	16.62 $\pm$ 0.64	49.61 $\pm$ 1.54	17.89 $\pm$ 0.42	52.77 $\pm$ 1.01	<b>9.96</b> $\pm$ 0.37	32.92 $\pm$ 1.05	<u>19.11</u> $\pm$ 0.40	<u>56.28</u> $\pm$ 1.00
DeepSurv	16.55 $\pm$ 0.93	49.85 $\pm$ 3.02	17.80 $\pm$ 0.49	52.62 $\pm$ 1.25	10.09 $\pm$ 0.38	33.28 $\pm$ 1.15	19.20 $\pm$ 0.41	56.48 $\pm$ 1.08
CoxTime	16.54 $\pm$ 0.83	49.67 $\pm$ 2.67	17.80 $\pm$ 0.58	52.56 $\pm$ 1.47	10.28 $\pm$ 0.45	34.18 $\pm$ 1.53	19.17 $\pm$ 0.40	56.45 $\pm$ 1.10
DeepHit	17.50 $\pm$ 0.83	52.10 $\pm$ 2.16	19.61 $\pm$ 0.38	56.67 $\pm$ 1.10	11.83 $\pm$ 0.39	37.72 $\pm$ 1.02	20.66 $\pm$ 0.32	60.06 $\pm$ 0.72
DeepEH	16.56 $\pm$ 0.65	49.42 $\pm$ 1.53	17.62 $\pm$ 0.52	<u>52.08</u> $\pm$ 1.27	10.11 $\pm$ 0.37	33.30 $\pm$ 1.10	19.30 $\pm$ 0.39	56.67 $\pm$ 0.94
SuMo-net	16.49 $\pm$ 0.83	49.74 $\pm$ 2.21	17.77 $\pm$ 0.47	52.62 $\pm$ 1.11	10.07 $\pm$ 0.40	33.20 $\pm$ 1.10	19.40 $\pm$ 0.38	56.87 $\pm$ 0.96
SODEN	16.52 $\pm$ 0.63	49.39 $\pm$ 1.97	<b>17.05</b> $\pm$ 0.63	<b>50.45</b> $\pm$ 1.97	10.13 $\pm$ 0.24	33.37 $\pm$ 0.57	19.07 $\pm$ 0.50	56.15 $\pm$ 1.35
SurvNode	16.67 $\pm$ 1.32	49.73 $\pm$ 3.89	17.42 $\pm$ 0.53	51.70 $\pm$ 1.16	10.40 $\pm$ 0.29	34.37 $\pm$ 1.03	19.58 $\pm$ 0.34	57.49 $\pm$ 0.84
DCM	16.58 $\pm$ 0.87	49.48 $\pm$ 2.23	17.66 $\pm$ 0.54	52.26 $\pm$ 1.23	10.13 $\pm$ 0.50	33.40 $\pm$ 1.38	19.29 $\pm$ 0.42	56.68 $\pm$ 1.09
DeSurv	16.71 $\pm$ 0.75	49.61 $\pm$ 2.15	17.98 $\pm$ 0.46	53.23 $\pm$ 1.15	10.06 $\pm$ 0.62	33.18 $\pm$ 1.93	19.50 $\pm$ 0.40	57.28 $\pm$ 0.89
<b>NFM-PF</b>	<u>16.33</u> $\pm$ 0.75	<u>49.07</u> $\pm$ 1.96	<u>17.60</u> $\pm$ 0.55	52.12 $\pm$ 1.34	<b>9.96</b> $\pm$ 0.39	<b>32.84</b> $\pm$ 1.15	19.14 $\pm$ 0.39	56.35 $\pm$ 1.00
<b>NFM-FN</b>	<b>16.11</b> $\pm$ 0.81	<b>48.21</b> $\pm$ 2.04	17.66 $\pm$ 0.52	52.41 $\pm$ 1.22	10.05 $\pm$ 0.39	33.11 $\pm$ 1.10	<b>18.97</b> $\pm$ 0.60	<b>55.87</b> $\pm$ 1.50

Empirical evaluations over 2 datasets with larger scale

Model	MIMIC-III		KKBOX	
	IBS	INBLL	IBS	INBLL
CoxPH	20.40 $\pm$ 0.00	60.02 $\pm$ 0.00	12.60 $\pm$ 0.00	39.40 $\pm$ 0.00
GBM	17.70 $\pm$ 0.00	52.30 $\pm$ 0.00	11.81 $\pm$ 0.00	38.15 $\pm$ 0.00
RSF	17.79 $\pm$ 0.19	53.34 $\pm$ 0.41	14.46 $\pm$ 0.00	44.39 $\pm$ 0.00
DeepSurv	18.58 $\pm$ 0.92	55.98 $\pm$ 2.43	11.31 $\pm$ 0.05	35.28 $\pm$ 0.15
CoxTime	17.68 $\pm$ 1.36	52.08 $\pm$ 3.06	10.70 $\pm$ 0.06	33.10 $\pm$ 0.21
DeepHit	19.80 $\pm$ 1.31	59.03 $\pm$ 4.20	16.00 $\pm$ 0.34	48.64 $\pm$ 1.04
SuMo-net	18.62 $\pm$ 1.23	54.51 $\pm$ 2.97	11.58 $\pm$ 0.11	36.61 $\pm$ 0.28
DCM	18.02 $\pm$ 0.49	52.83 $\pm$ 0.94	10.71 $\pm$ 0.11	33.24 $\pm$ 0.06
DeSurv	18.19 $\pm$ 0.65	54.69 $\pm$ 2.83	10.77 $\pm$ 0.21	33.22 $\pm$ 0.10
<b>NFM-PF</b>	<b>16.28</b> $\pm$ 0.36	<b>49.18</b> $\pm$ 0.92	11.02 $\pm$ 0.11	35.10 $\pm$ 0.22
<b>NFM-FN</b>	<b>17.47</b> $\pm$ 0.45	<b>51.48</b> $\pm$ 1.23	<b>10.63</b> $\pm$ 0.08	<b>32.81</b> $\pm$ 0.14



- We have introduced NFM as a flexible and powerful neural modeling framework for survival analysis.
- NFM is shown to be both statistically correct in theory, and empirically effective in predictive tasks.
- Future directions: Establishing theory guarantees toward more realistic predictive metrics instead of nonparametric parameter estimation.



Steve Bennett.

Analysis of survival data by the proportional odds model.

*Statistics in medicine*, 2(2):273–277, 1983.



David R Cox.

Regression models and life-tables.

*Journal of the Royal Statistical Society: Series B (Methodological)*,  
34(2):187–202, 1972.



David Faraggi and Richard Simon.

A neural network model for survival data.

*Statistics in medicine*, 14(1):73–82, 1995.



Philip Hougaard.

Survival models for heterogeneous populations derived from stable distributions.

*Biometrika*, 73(2):387–396, 1986.



Jared L Katzman, Uri Shaham, Alexander Cloninger, Jonathan Bates,  
Tingting Jiang, and Yuval Kluger.

Deepsurv: personalized treatment recommender system using a cox proportional hazards deep neural network.

*BMC medical research methodology*, 18(1):1–12, 2018.