

Interpretable Deep Survival Learning for Credit Default Risk Prediction in Financial Markets

Elliot C. Love

Department of Computer Science, University of New Hampshire, Durham, NH, USA.
love38@unh.edu

Junguo Peng

Department of Computer Science, Colorado State University, Fort Collins, CO, USA.
junguop@colostate.edu

Claude Murray

Department of Computer Science, University of Alabama at Birmingham, Birmingham, AL, USA.
contactclaude@uab.edu

Yash Kingh

Department of Electrical Engineering and Computer Science, University of Missouri, Columbia, MO, USA.
yashwork@missouri.edu

Abstract

Credit default risk prediction remains a central challenge in financial markets, where the temporal dynamics of borrower behavior and the need for regulatory transparency demand models that are both accurate and interpretable. Traditional survival analysis methods offer a principled way to model time-to-event data but often fail to capture complex nonlinear interactions in high-dimensional financial datasets. Deep learning approaches, particularly neural networks adapted for survival analysis, have demonstrated superior predictive performance; however, their black-box nature undermines trust and compliance with emerging explainability mandates. This paper presents a comprehensive framework for interpretable deep survival learning applied to credit default risk, integrating architectural innovations with system-level deployment considerations. We examine the trade-offs between model complexity and interpretability, discussing attention mechanisms, time-dependent gradient-based explanations, and surrogate modeling as pathways to transparency. The discussion extends to infrastructure challenges, including real-time inference pipelines, data governance, model sustainability, and fairness across demographic segments. We argue that interpretability is not merely a technical add-on but a structural requirement for robust risk management, regulatory alignment, and equitable credit access. By synthesizing insights from survival analysis, deep learning, and socio-technical systems, we propose a multi-layered evaluation framework that balances predictive power with actionable explanations. The paper further explores policy implications, particularly in light of evolving regulations such as the European Union's AI Act, and highlights the role of interpretable survival models in reducing systemic risk. Through case illustrations and cross-domain comparisons, we demonstrate that interpretable deep survival learning can bridge the gap between high-performance forecasting and responsible financial governance.

Keywords

interpretable machine learning, survival analysis, credit default risk, deep learning, financial system governance, fairness, model deployment.

1. Introduction

The accurate prediction of credit default risk is fundamental to the stability of financial markets and the efficient allocation of capital. Lending institutions, regulatory bodies, and investors rely on risk assessment models to determine the likelihood that a borrower will fail to meet their obligations over a given time horizon. Traditional approaches, such as logistic regression and Cox proportional hazards models, have long served as the backbone of credit scoring systems because of their statistical interpretability and regulatory acceptance [1]. However, these methods impose strong parametric assumptions and often underperform in the presence of high-dimensional, heterogeneous, and dynamically evolving financial data. The advent of deep learning has provided a pathway to capture nonlinear relationships, temporal dependencies, and complex interactions among borrower attributes, macroeconomic indicators, and market conditions [2]. Yet the opacity of deep neural networks poses serious challenges for model justification, auditability, and fairness, especially in regulated financial environments where explainability is increasingly mandated.

Survival analysis offers a natural framework for credit default prediction because it models the distribution of time until an event of interest, such as default, while accounting for censored observations where the event has not occurred by the end of the observation period. Deep survival learning extends this framework by embedding survival objectives into neural architectures, enabling the learning of flexible hazard functions from raw data [3]. Models such as DeepSurv, Cox-Time, and neural multi-task logistic regression have shown substantial improvements in concordance indices and calibration over classical benchmarks [4]. Nevertheless, the interpretability of these models remains an open problem. Recent work [5] has compared transformer-based architectures with classical time-series models for financial forecasting, highlighting that while transformers achieve higher accuracy, their complex attention mechanisms introduce additional layers of opacity. This trade-off between predictive power and transparency is at the core of our inquiry.

The objective of this paper is to examine interpretable deep survival learning as a socio-technical system for credit default risk prediction, moving beyond isolated model accuracy to consider architectural design choices, deployment infrastructure, regulatory compliance, and fairness implications. We argue that interpretability must be designed into the system from the ground up, rather than retrofitted as an afterthought. To this end, we explore a range of interpretability techniques, including attention-based feature attribution, time-dependent explainability via gradient integration, and the use of inherently interpretable surrogate models. We also address the practical challenges of deploying such models in real-time risk management pipelines, including latency constraints, data versioning, and monitoring for concept drift. By situating our analysis within the broader context of financial market governance, we highlight how interpretable survival models can contribute to systemic stability and equitable access to credit.

The remainder of the paper is organized as follows. Section 2 reviews the background of survival analysis in credit risk and recent advances in interpretable machine learning. Section 3 presents a detailed architectural framework for interpretable deep survival learning, focusing on model components and explanation mechanisms. Section 4 discusses

infrastructure and deployment considerations, including data pipelines, model sustainability, and operational trade-offs. Section 5 addresses robustness, fairness, and policy implications, drawing on regulatory frameworks and ethical principles. Section 6 concludes with a synthesis of findings and directions for future research.

2. Background and Related Work

Credit default risk modeling has evolved from simple linear discriminant analysis to sophisticated machine learning ensembles. Survival analysis contributes a distinctive perspective by treating default as a time-to-event outcome, allowing the estimation of survival probabilities over continuous time horizons [6]. The Cox proportional hazards model remains widely used due to its semi-parametric nature and interpretable hazard ratios. However, its proportional hazards assumption is often violated in practice, leading to biased estimates when covariates effects change over time [7]. Deep learning methods relax these assumptions by allowing non-linear hazard functions and time-varying effects through neural network parameterizations.

DeepSurv, introduced by Katzman et al. [8], combined a Cox partial likelihood loss with a feedforward neural network, demonstrating improved predictive performance on survival datasets. Subsequent work has extended this approach to handle competing risks, time-dependent covariates, and recurrent events. More recently, attention-based architectures have been applied to survival analysis, leveraging the ability of transformers to model long-range dependencies in sequential borrower data [9]. Despite these advances, the interpretability of deep survival models has lagged behind. Post-hoc explanation methods, such as SHAP and LIME, have been adapted to survival settings, but they often produce inconsistent explanations across time points and fail to capture the temporal dynamics of risk factors [10].

Interpretable machine learning has emerged as a distinct subfield, driven by regulatory demands and societal concerns. The European Union’s General Data Protection Regulation established a right to explanation for automated decisions, while the proposed AI Act introduces stricter transparency requirements for high-risk AI systems, including credit scoring [11]. In the financial domain, supervisory authorities such as the Federal Reserve and the European Banking Authority have emphasized the need for model interpretability to ensure fair lending practices and prudent risk management. Several lines of research have pursued intrinsically interpretable models, such as generalized additive models with pairwise interactions (GA2Ms) and neural additive models, which achieve competitive accuracy while maintaining feature-level transparency [12]. However, these models often require careful feature engineering and may not capture the full complexity of survival dynamics.

A parallel stream of work has focused on time-dependent interpretability, recognizing that the importance of a borrower’s attributes may change over the prediction horizon. Qin et al. [13] proposed an interpretable deep survival analysis framework that integrates extreme gradient boosting with time-dependent feature importance, allowing analysts to trace how each variable contributes to the survival function at different time points. This approach is particularly relevant for credit default, where early warning signals may differ from long-term risk drivers. Similarly, attention mechanisms in transformer-based survival models can provide temporal attribution scores, highlighting which historical events most influence the current hazard [14]. Nevertheless, these methods raise questions about the stability and faithfulness of explanations, especially when model confidence is low or when data distributions shift.

In the broader landscape of financial machine learning, reinforcement learning has been applied to dynamic portfolio optimization, where attention mechanisms enhance the model’s ability to focus on relevant market signals [15]. While portfolio optimization and credit risk share some methodological foundations, the interpretability requirements differ: credit decisions directly affect individuals’ access to financial services and are subject to anti-discrimination laws. Thus, interpretability in survival-based credit models must go beyond global feature importance to provide local, individualized explanations that can be audited for fairness.

3. Interpretable Deep Survival Learning: Architectural Framework

An interpretable deep survival learning system for credit default prediction must integrate three core capabilities: flexible hazard modeling, transparent feature attribution, and temporally coherent explanation. We propose a modular framework that decomposes the prediction pipeline into four components: data preprocessing, a deep survival encoder, an explanation module, and a decision layer. The encoder can be implemented using a variety of neural architectures, including fully connected networks, recurrent networks for sequential data, and transformers for capturing long-term dependencies. Each architecture imposes different trade-offs between capacity and interpretability. For instance, recurrent networks naturally model temporal sequences but are difficult to explain beyond attention over time steps, while transformers provide explicit attention weights over input tokens but may yield diffuse attention patterns when input sequences are noisy [16].

The explanation module is the critical new component that distinguishes an interpretable system from a black-box one. We advocate for a hybrid approach that combines post-hoc and intrinsic interpretability methods. For the encoder, we can incorporate attention layers that output temporal and feature-level attention weights. These weights serve as a first-order explanation of which inputs the model considers important at each prediction step. However, attention is not always faithful to the model’s true decision process, and multiple studies have shown that attention weights can be manipulated without affecting predictions [17]. Therefore, we supplement attention with gradient-based sensitivity analysis, such as integrated gradients, applied to the survival loss function. By integrating the gradient of the predicted hazard with respect to input features along a path from a baseline, we obtain a time-dependent attribution score that satisfies axioms of completeness and sensitivity [18].

For scenarios where regulators demand a fully transparent decision rule, we propose using a distilled surrogate model. After training the deep survival encoder, we fit an interpretable model, such as a Cox proportional hazards model with selected interactions or a risk-calibrated piecewise exponential model, to approximate the encoder’s predictions in the neighborhood of each instance. The surrogate model provides a locally faithful explanation that can be inspected by human analysts. To ensure that the surrogate does not systematically deviate from the deep model, we monitor the approximation error using divergence metrics and retrain the surrogate if drift is detected [19].

The decision layer integrates the survival probabilities and their explanations into a deployable scoring system. Instead of outputting a single default probability, the framework produces a survival curve along with a ranked list of risk factors and their temporally varying contributions. This output format aligns with regulatory expectations for model transparency and allows loan officers to understand why a particular borrower was assigned a high risk score. Moreover, the framework supports counterfactual explanations: by slightly perturbing

input features, one can generate “what-if” scenarios that show how changes in a borrower’s attributes would affect their predicted survival probability [20].

4. Infrastructure, Deployment, and System-Level Considerations

Deploying an interpretable deep survival learning model in a production credit risk environment introduces infrastructural challenges that are often underestimated in academic research. Real-time risk assessment requires low-latency inference, typically under 100 milliseconds for point-of-sale credit decisions. Deep neural networks, especially transformer-based models, impose significant computational overhead. To meet latency constraints while preserving accuracy, we consider model compression techniques such as pruning, quantization, and knowledge distillation. Distillation is particularly attractive because it can produce a smaller, faster model that retains the interpretability characteristics of the teacher network [21]. However, care must be taken to ensure that compression does not amplify biases or degrade explanation faithfulness.

Data governance is another critical dimension. Credit default models ingest a wide array of sensitive data, including demographic attributes, transaction histories, and alternative data sources such as utility payments. The interpretability framework must be integrated with data lineage tracking to ensure that explanations can be traced back to the specific features and preprocessing steps that contributed to a prediction. Version control for both data and model artifacts is essential for auditability, especially when regulators request retrospective explanations for past decisions [22]. Moreover, the system must handle censored data correctly during training and inference. In a dynamic deployment, new borrowers are observed over time, and their censoring status changes as the prediction horizon extends. The model must be updated regularly, either through periodic retraining or online learning, to reflect evolving borrower populations and macroeconomic conditions.

Sustainability of the model over long time horizons is a concern. Financial markets experience structural breaks, regulatory changes, and unexpected shocks (e.g., a pandemic) that can render previously learned patterns obsolete. The interpretability framework can aid in detecting concept drift by monitoring the distribution of feature attributions. If the importance of a previously stable risk factor suddenly shifts, it signals a need for model recalibration. Building automated drift detection and alerting mechanisms into the deployment pipeline reduces the risk of silent model degradation [23].

Cross-domain comparisons reveal that credit risk systems share infrastructural patterns with other socio-technical domains, such as healthcare and predictive maintenance. In healthcare, interpretable survival models are used to estimate patient prognosis while providing clinicians with understandable risk factors. The evaluation metrics and governance frameworks from healthcare, such as the requirement for prospective validation and fairness audits, offer valuable lessons for financial applications. However, financial systems operate under tighter latency constraints and face stronger adversarial pressure from borrowers attempting to manipulate their risk profiles. Therefore, adversarial robustness testing must be integrated into the deployment lifecycle, examining how sensitive the model’s explanations are to small perturbations in input features [24].

5. Robustness, Fairness, and Policy Implications

Interpretable deep survival learning for credit default must be robust to distributional shifts and adversarial perturbations to maintain trustworthiness over time. A model that is highly accurate on historical data but fails under stress scenarios can exacerbate systemic risk.

Robustness can be enhanced through adversarial training, where the model is trained on perturbed examples that simulate plausible changes in borrower behavior. Interpretability aids robustness by providing a diagnostic tool: if explanations become unstable under slight input variations, it indicates that the model's decision boundary is fragile [25]. Regular stress testing of explanations, similar to stress testing of risk models, should become a standard practice.

Fairness is a paramount concern in credit risk modeling. Survival models that are trained on historical data may inherit biases present in the lending process, such as disparate denial rates across racial or gender groups. Even if demographic attributes are excluded from the model, proxy variables like zip code or income can encode historical discrimination. Interpretable survival models allow regulators and auditors to examine the extent to which each feature contributes to risk scores for different demographic groups. Techniques such as fairness-aware survival objectives, which penalize differences in survival calibration across groups, can be integrated into the training process [26]. The temporal nature of survival analysis adds complexity: a model that is fair at a one-year horizon may be unfair at a five-year horizon if hazard functions diverge. Therefore, fairness evaluation must be performed at multiple prediction windows, and the interpretability framework should enable disaggregated analysis of survival curves by protected attributes.

Policy implications are far-reaching. Regulatory frameworks such as the European Union's AI Act classify credit scoring systems as high-risk AI, requiring transparency, human oversight, and the right to explanation. Interpretable deep survival models can help financial institutions meet these requirements without sacrificing predictive performance. However, there is a risk that regulators may demand explanations that are causally valid, whereas most post-hoc methods only provide associational attribution. Bridging the gap between associational and causal explanations remains an open research challenge. One promising direction is the use of counterfactual explanations derived from structural causal models of borrower behavior [27].

From a governance perspective, financial institutions must establish internal committees that oversee the development, validation, and monitoring of interpretable survival models. These committees should include data scientists, domain experts, ethicists, and legal professionals. The interpretability framework should produce documentation that is accessible to non-technical stakeholders, including loan officers and regulatory examiners. The adoption of model cards and datasheets, as advocated by [28], provides a structured format for communicating model purpose, performance, limitations, and fairness evaluations.

6. Conclusion

Interpretable deep survival learning represents a significant advancement for credit default risk prediction, offering the ability to model complex, time-dependent dynamics while providing transparency into the decision-making process. This paper has argued that interpretability must be treated as an integral design principle rather than a secondary consideration, requiring systematic attention to architectural choices, deployment infrastructure, fairness constraints, and regulatory compliance. The proposed framework, combining attention mechanisms, time-dependent gradient integration, and surrogate distillation, provides a practical pathway for deploying high-performance survival models in real-world financial systems. However, open challenges remain, including the causal validation of explanations, the stability of attributions under distribution shift, and the scalability of interpretability methods to large-scale portfolios. Future research should focus

on developing evaluation metrics that jointly measure predictive accuracy and explanation quality, as well as on exploring federated learning paradigms that preserve interpretability across distributed data silos. As financial markets continue to digitize and regulatory expectations tighten, interpretable deep survival learning will be essential for building credit systems that are both powerful and responsible.

References

1. Katzman, J. L., Shaham, U., Cloninger, A., Bates, J., Jiang, T., & Kluger, Y. (2018). DeepSurv: Personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC Medical Research Methodology*, 18(1), 24.
2. Cox, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2), 187–202.
3. Lee, C., Zame, W. R., Yoon, J., & van der Schaar, M. (2018). DeepHit: A deep learning approach to survival analysis with competing risks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).
4. Kvamme, H., Borgan, O., & Scheel, I. (2019). Time-to-event prediction with neural networks and Cox regression. *Journal of Machine Learning Research*, 20(129), 1–30.
5. Liu, T. (2026). A Comparative Study of Transformer-Based and Classical Models for Financial Time-Series Forecasting. *Journal of Risk and Financial Management*, 19(3), 203.
6. Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765–4774.
7. European Commission. (2021). Proposal for a regulation laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). COM(2021) 206 final.
8. Agarwal, R., Melnick, L., Frosst, N., Zhang, X., Lengerich, B., Caruana, R., & Hinton, G. E. (2021). Neural additive models: Interpretable machine learning with neural nets. *Advances in Neural Information Processing Systems*, 34, 4699–4711.
9. Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. *International Conference on Learning Representations*.
10. Qin, X., Yu, R., Khayati, A., Qiu, Z., Zou, G., Li, Y., & Wang, L. (2025, November). Interpretable and Interactive Deep Survival Analysis with Time-dependent EXtreme Gradient Integration. In *2025 IEEE International Conference on Data Mining (ICDM)* (pp. 673-682). IEEE.
11. Sundararajan, M., Taly, A., & Yan, Q. (2017). Axiomatic attribution for deep networks. *Proceedings of the 34th International Conference on Machine Learning*, 70, 3319–3328.
12. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?”: Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.
13. Liu, T. (2026). Interpretable Machine Learning for Volatility Forecasting Under Realistic Walk-Forward Constraints.
14. Xue, P., & Ye, Y. (2026). Attention-enhanced reinforcement learning for dynamic portfolio optimization. *Intelligent Systems with Applications*, 200622.

15. Liu, T. (2022, December). Financial Constraint'Impact on Firms' ESG Rating Based on Chinese Stock Market. In 2022 4th International Conference on Economic Management and Cultural Industry (ICEMCI 2022) (pp. 1085-1095). Atlantis Press.
16. Liu, T. (2026). Volatility Forecasting and Early-Warning Market Stress Detection: A Leakage-Safe Evaluation with Tree Ensembles and Transformers.
17. Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., & Elhadad, N. (2015). Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1721–1730.
18. Horel, J. A., & Geman, S. (2022). A computational framework for model stability. Proceedings of the National Academy of Sciences, 119(20), e2120967119.
19. Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608.
20. Selbst, A. D., Boyd, D., Friedler, S. A., Venkatasubramanian, S., & Vertesi, J. (2019). Fairness and abstraction in sociotechnical systems. Proceedings of the Conference on Fairness, Accountability, and Transparency, 59–68.
21. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... & Gebru, T. (2019). Model cards for model reporting. Proceedings of the Conference on Fairness, Accountability, and Transparency, 220–229.
22. Pearl, J. (2009). Causality: Models, reasoning, and inference (2nd ed.). Cambridge University Press.
23. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785–794.