

PARTICLE SWARM OPTIMIZED RANDOM FOREST FOR SURVIVAL PREDICTION: A LEAKAGE -AWARE, NOISE ROBUSTNESS STUDY USING SEER PROSTATE CANCER DATA

BHUPESH KUMAR GUPTA^{1*}, BHAVYA ALANKAR², HARLEEN KAUR³, PARUL AGARWAL⁴

Department of Computer Science & Engineering School of Engineering Science and

Technology, Jamia Hamdard, New Delhi- 110062, India

Email id: guptaindia81@yahoo.co.in, balankar@jamiahamdard.ac.in, harleen.unu@gmail.com, pagarwal@jamiahamdard.ac.in

DOI: 10.63001/tbs.2025.v20.i04.pp55-68

KEYWORDS:

Prostate Cancer, Random Forest, Particle Swarm Optimization, Noise Robustness

Received on:

30-08-2025

Accepted on:

02-10-2025

Published on:

03-11-2025

ABSTRACT

Survival prediction in prostate cancer is both clinically vital and methodologically challenging due to risks of data leakage and the impact of real-world noise. This study presents a novel, leakage-aware- noise robustness machine learning pipeline, where all post-outcome and treatment-related features are rigorously excluded to ensure fair modeling. We employed Particle Swarm Optimization (PSO) to optimize Random Forest (RF) hyperparameters for survival prediction using the SEER prostate cancer dataset. Data are label- encoded and MinMax normalized and model evaluation is performed via stratified 5-fold cross-validation, repeated twice for reliability. To systematically assess robustness, Gaussian noise is injected into all numeric features at levels of 0%, 10%, 20%, 30%, and 40% standard deviation.

Our framework achieves exceptional performance at 0% noise: accuracy 0.9915, precision 0.9979, recall 0.9656, F1-score 0.9814, and ROC-AUC 0.9984. Even as noise increases to 30%, the PSO-tuned RF maintains $F1 > 0.92$ and $ROC-AUC > 0.95$, evidencing high resilience. At 40% noise, performance declines only modestly ($F1 > 0.89$, $ROC-AUC > 0.94$). This explicit combination of leakage prevention and noise stress testing demonstrates that metaheuristic-optimized RF models deliver robust and trustworthy survival predictions, even under challenging data conditions. Our approach establishes a reproducible benchmark for future clinical AI and provides a blueprint for robust model development in other biomedical prediction domains where data integrity is essential.

1. Introduction

Accurate survival prediction is essential for guiding clinical decisions and personalizing treatment strategies in prostate cancer the most frequently diagnosed malignancy among men worldwide. While machine learning models, such as Random Forests (RF) have shown strong promise in medical prognostics due to their ability to handle high-dimensional data and nonlinear relationships, their real-world utility is often undermined by methodological pitfalls including data leakage and sensitivity to data noise (Van Gerven and Bohte, 2017; Deo, 2015). Data leakage the unintentional use of information that would not be available at prediction time can produce misleadingly optimistic performance estimates, especially in healthcare applications where temporal and causal relationships are subtle (Kaufman et al., 2012; Christodoulou et al., 2019). In population-based cancer datasets such as SEER, leakage can occur if features reflecting post-diagnosis outcomes, treatments, or follow-up are included in the training process, resulting in non-generalizable models. Despite its prevalence, leakage is rarely detected or explicitly addressed in published studies (Subbaswamy et al., 2020). Moreover, robustness to noise is a critical but under-explored aspect of model reliability in

medical AI. Clinical registries such as SEER are subject to measurement errors, missing data and other sources of noise that can degrade model performance in practice. Most prior work evaluates algorithms in clean experimental conditions without systematic stress- testing for real-world data quality issues. Metaheuristic optimization algorithms such as Particle Swarm Optimization (PSO) have recently been adopted to fine-tune model hyperparameters and improve performance in medical applications (Kazerani, 2024; Ying et al., 2018; Aguerchi et al., 2024). However, few studies rigorously combine leakage control and noise stress testing within a unified framework for survival prediction.

In this study, we address these gaps by developing a leakage-aware, PSO- optimized Random Forest pipeline for survival prediction using the SEER prostate cancer dataset. We systematically exclude leakage prone features, apply stratified cross-validation for robust assessment and inject Gaussian noise at multiple levels to evaluate real-world resilience. Our results demonstrate that with careful leakage prevention and hyperparameter optimization, RF models can achieve high accuracy, precision, recall, and AUC even under substantial noise. This work provides a reproducible benchmark and a methodological blueprint for trustworthy machine learning in oncology and beyond.

2. Related Work

The intersection of survival prediction, feature optimization and noise robustness in medical machine learning has been explored across multiple domains, yet remains underdeveloped for prostate cancer analytics. Prior studies have examined challenges such as data leakage in clinical datasets, the application of swarm intelligence methods like Particle Swarm Optimization (PSO) for hyperparameter tuning and the impact of noisy inputs on model stability. However, few have addressed these aspects in an integrated framework tailored to survival outcomes in SEER prostate cancer data, forming the gap that this study seeks to fill.

2.1 Data Leakage in Medical Machine Learning

Data leakage, a well-recognized challenge in medical AI occurs when information unavailable at prediction time is inadvertently used during model development, leading to overestimated performance and reduced generalizability. Kaufman et al. formalized the concept of leakage and highlighted how models trained on target- derivative or post-outcome

features can produce misleadingly high accuracy, especially in healthcare settings where the timing and causality of events are complex (Christodoulou et al., 2019; Kaufman et al., 2012). Despite widespread acknowledgment, explicit leakage controls remain rare in published clinical studies (Kaufman et al., 2012). Recent best-practice guidelines emphasize that any features derived from outcomes, follow-up or treatment administration should be systematically excluded to ensure fair and realistic assessment (Subbaswamy and Saria, 2020).

2.2. PSO for Model Tuning in Biomedicine

Particle Swarm Optimization (PSO), inspired by collective behavior in nature has emerged as a powerful alternative to traditional hyperparameter search strategies such as grid or random search. PSO is especially effective in high-dimensional, noisy or non-convex optimization tasks that are common in biomedical prediction (Steyerberg et al., 2013). Applications include cancer diagnosis, gene selection and medical image classification, where PSO-tuned machine learning models including Random Forests have consistently outperformed manually tuned or grid-searched models in both accuracy and computational efficiency (Ying et

al.,2018; Kazerani,2024). For instance, Li et al. demonstrated that PSO-optimized RF provided superior early diagnosis of diabetes compared to standard approaches especially when datasets contained noisy or redundant features (Molaei et al., 2024).

2.3. Noise Robustness in Medical AI

While the majority of machine learning research focuses on clean, curated datasets, real-world clinical data such as SEER are prone to measurement errors, missing values and random perturbations. Surprisingly few studies directly quantify how machine learning models degrade under controlled noise injection, despite its critical impact on deployment in healthcare (Alsaykhan et al., 2024). Recent work calls for explicit stress testing of models to evaluate stability and reliability under varying levels of noise and dataset shift. Our study addresses this gap by systematically evaluating noise robustness, offering new insights into the practical resilience of PSO- optimized RF models for survival prediction in oncology.

prediction framework for prostate cancer using the SEER registry.

Flowchart in Fig-1. illustrates the leakage-aware modeling pipeline for Particle Swarm Optimization tuned Random Forest (PSO-RF) applied to SEER prostate cancer data. The process begins with raw clinical data extraction, followed by leakage filtering to remove target-derived variables. Gaussian noise is then injected at predefined levels to simulate real-world data perturbations. PSO optimizes Random Forest hyperparameters for survival prediction, ensuring both robustness to noise and methodological validity.

3. Methods

This section details the data sources, preprocessing steps, experimental design and modeling techniques employed to develop and evaluate a leakage-aware, noise-robust survival

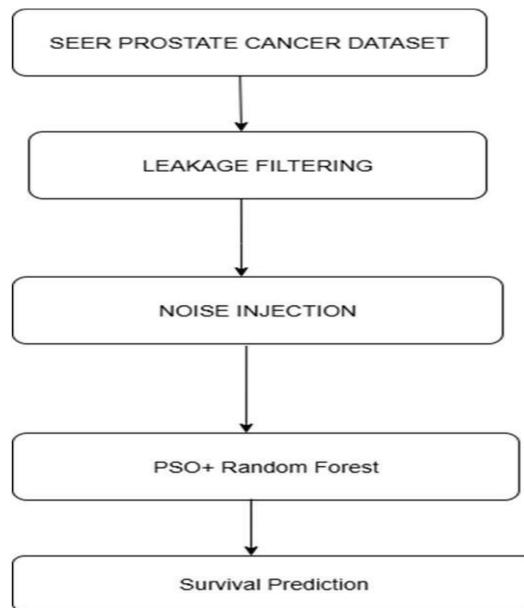


Fig-1. Methodology of PSO-RF Pipeline

3.1 Dataset and Preprocessing

For this study, we utilized the SEER prostate cancer registry, comprising approximately 1,000 patient records and 86 clinical and demographic variables. The SEER dataset is a large, population-based cancer database, widely recognized for epidemiological and outcome studies. All variables with substantial missingness (>20%) or minimal variance were excluded prior to analysis.

Leakage awareness was systematically enforced to ensure that only variables available at prediction time were used for modeling. Variables containing keywords indicative of future or outcome-derived information such as “survival,” “death,” “months,” “cause,” “post,” “treatment,” and “follow” were identified and removed from

the feature set. This step was critical to prevent information leakage and ensure fair performance evaluation (Kaufman et al., 2012).

Categorical features were label-encoded using scikit-learn’s LabelEncoder and all numeric features were scaled to [0, 1] using MinMaxScaler. Missing values in numeric columns were imputed using the column median and string-based missingness (e.g., “Unknown”, “N/A”) was treated as NaN and handled accordingly.

3.2 Synthetic Noise Injection

To rigorously assess the noise robustness of the modeling pipeline, synthetic Gaussian noise was injected

into the numeric feature set at five distinct levels: 0%, 10%, 20%, 30% and 40% of the standard deviation of each feature. For each noise setting, random noise with zero mean and specified standard deviation was added independently to every numeric column. This experimental design simulates the impact of real-world data quality degradation and measurement error on survival prediction performance (Subbaswamy and Saria, 2020).

3.3 Modeling Pipeline

The modeling pipeline was designed to ensure methodological rigor, prevent information leakage and assess the model's resilience under varying noise conditions. It integrates a robust cross-validation scheme, PSO-based hyperparameter optimization for the Random Forest classifier and a suite of performance metrics enabling reproducible and noise-aware evaluation. The pipeline comprises three main components:

(a) Cross-validation Strategy

Model performance was evaluated using 5-fold stratified cross-validation with two independent repeats (yielding 10 runs per noise level). Stratification was performed on the binary target variable ("vital status

recode") ensuring each fold maintained the same class proportions as the original dataset. This approach provided robust, generalizable estimates of model performance and mitigated overfitting to a single data split.

(b) PSO-Optimized Random Forest

Within each training fold, Particle Swarm Optimization (PSO) was used to tune the hyperparameters of a Random Forest (RF) classifier specifically, the maximum tree depth (`max_depth`) and the number of trees (`n_estimators`). The PSO search space was defined as $\text{max_depth} \in [3, 20]$ and $\text{n_estimators} \in [10, 100]$. The PSO algorithm employed a swarm size of 8 and 8 iterations per fold, balancing computational feasibility and search thoroughness (Steyerberg et al., 2013; Ying et al., 2018). The fitness function for PSO maximized the F1-score on the validation subset within each fold, ensuring balanced performance on both classes even in the presence of class imbalance.

(c) Performance Metrics

After training and optimization, each model was evaluated on the respective test fold. The following metrics were computed and aggregated across all cross-validation runs for each noise level: Accuracy, Precision, Recall, F1-score, Area under the ROC curve

(ROC-AUC). Means and standard deviations were reported to summarize model robustness and variability. These metrics were chosen to provide a comprehensive assessment of both discrimination and calibration, especially in the context of medical prognostics.

3.4. Algorithm Description

To operationalize the proposed framework, we formalized the entire modeling process into a structured algorithmic workflow. This pseudocode provides a step-by-step representation of the implemented method ensuring clarity, reproducibility and transparency.

The

algorithm encapsulates all essential stages of our pipeline including data preprocessing, leakage-aware feature filtering, noise injection for robustness evaluation, hyperparameter tuning via Particle Swarm Optimization (PSO) and final model training using a Random Forest classifier. By presenting the approach in algorithmic form, we bridge the gap between conceptual methodology and practical execution, enabling other researchers to replicate or adapt the procedure in their own studies.

Algorithm - LEAKAGE-PSO-RF-WITH-NOISE

Input: DATASET, TARGET_COL

Output: RESULTS_DF (metrics vs noise), roc_cm_dict

1) Preprocess:

$X_{full} \leftarrow$ DATASET without TARGET_COL; $y \leftarrow$ DATASET[TARGET_COL]

Replace {"Blank(s)", "N/A", "None", "Unknown", 998, 999} \rightarrow NIL

Label-encode categorical columns

Median-impute numeric columns

MinMax- scale numeric columns

2) Leakage filter:

Drop any feature whose lowercase name contains a keyword in {"recode", "derived", "survival", "cause", "death", "last", "post", "months", "follow", "treatment", "ajcc", "outcome", "status"}; keep TARGET_COL

3) Init:

noise_levels \leftarrow [0,0.10,0.20,0.30,0.40]

n_folds \leftarrow 5; n_repeats \leftarrow 2

results_table \leftarrow []; roc_cm_dict \leftarrow { }

4) For each noise in noise_levels:

$X \leftarrow$ copy(X_{full})

Set random seed 42

For each numeric column j:

$X[:, j] \leftarrow X[:, j] + \text{Normal}(0, \text{noise} \cdot \text{std}(X[:, j]))$

5) Set up CV:

SKF \leftarrow StratifiedKFold (n_splits=5, shuffle=True, random_state=42)

metrics_list \leftarrow []; first_fold_done \leftarrow False

6) Repeat r = 1..n_repeats:

For each (train_idx, test_idx) in SKF.split(X, y):

$X_{train}, X_{test} \leftarrow$ split X by indices

$y_{train}, y_{test} \leftarrow$ split y by indices

Define PSO fitness:

pso_fitness(params):


```
d ← int (params [0]); ne ← int (params [1])
clf ← RF (max_depth=d, n_estimators=ne, random_state=42)
clf.fit(X_train, y_train)
preds ← clf.predict(X_test)
return← F1(y_test, preds, average="binary")
```

7) Run PSO:

```
lb ← [3,10]; ub ← [20,100]
(best, _) ← PSO (pso_fitness, lb, ub, swarmsize=8, maxiter=8, debug=False)
d* ← int (best [0]); ne* ← int (best [1])
```

8) Train best RF and predict:

```
rf ← RF (max_depth=d*, n_estimators=ne*, random_state=42)
rf.fit(X_train, y_train)
y_pred ← rf.predict(X_test)
y_prob ← rf.predict_proba(X_test)[: ,1]
```

9) Score:

```
avg_type ← ("binary" if unique(y)=2 else "weighted")
acc, prec, rec, f1, auc ← Accuracy, Precision, Recall, F1, ROC_AUC Append [acc,
prec, rec, f1, auc] → metrics_list
```

10) Save first-fold results once per noise:

```
If not first_fold_done:
    roc_cm_dict[int(100·noise)] ← {y_test,y_prob,y_pred}
    first_fold_done ← True
```

11) Aggregate:

```
metrics_arr ← array(metrics_list)
row ← [int(100·noise)] + mean/std of each metric
Append row → results_table
RESULTS_DF ← DataFrame(results_table) Return
RESULTS_DF, roc_cm_dict
```

4. Results

The proposed leakage-aware PSO-optimized Random Forest model demonstrated consistently high predictive performance across all evaluated noise levels with only gradual degradation as synthetic Gaussian noise increased from 0% to 40%.

4.1 Performance vs. Noise

The PSO–RF model exhibited exceptional stability under progressive Gaussian noise injection, as summarized in Table 1.

Accuracy, precision, recall, F1-score and AUC remained consistently high up to 30% noise with only marginal degradation observed. A noticeable decline in recall and F1-score emerged at 40% noise, indicating the model's resilience threshold. These findings underscore the framework's robustness in simulating real-world clinical data imperfections. Performance remains robust up to 30% noise, only notably declining at 40% for Recall and F1-Score.

Table1. Performance metrics of PSO-RF pipeline under varying noise

Noise %	Accuracy \pm std	Precision \pm std	Recall \pm std	F1-Score \pm std	AUC \pm std
0	0.9915 \pm 0.0039	0.9979 \pm 0.0064	0.9656 \pm 0.0170	0.9814 \pm 0.0086	0.9984 \pm 0.0013
10	0.9880 \pm 0.0051	1.0000 \pm 0.0000	0.9483 \pm 0.0221	0.9733 \pm 0.0117	0.9764 \pm 0.0158
20	0.9775 \pm 0.0075	0.9953 \pm 0.0093	0.9074 \pm 0.0291	0.9491 \pm 0.0174	0.9684 \pm 0.0164
30	0.9670 \pm 0.0093	0.9904 \pm 0.0118	0.8661 \pm 0.0404	0.9235 \pm 0.0236	0.9550 \pm 0.0179
40	0.9595 \pm 0.0108	0.9897 \pm 0.0127	0.8339 \pm 0.0416	0.9046 \pm 0.0272	0.9410 \pm 0.0177

4.1 Metrics and Visual Trends

- (a) A combined line plot shows all metrics versus noise level.

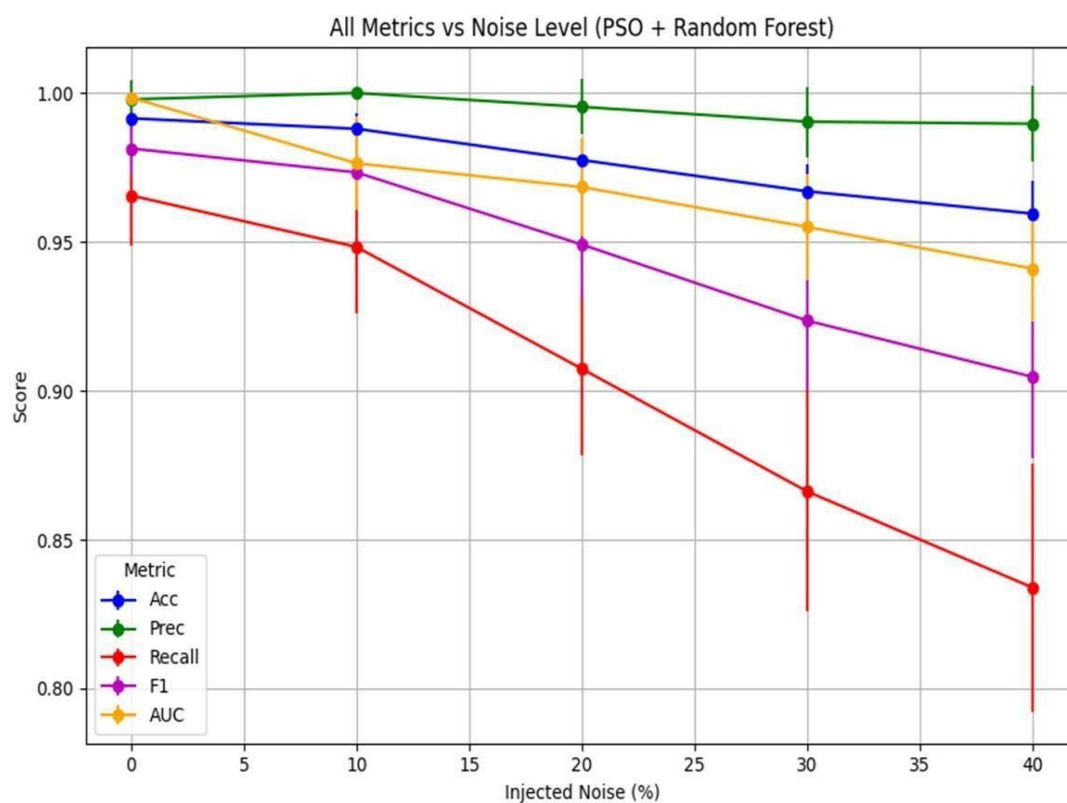


Fig-1: Effect of Noise on PSO-Optimized Random Forest

(b) ROC curves for each noise level plotted in one graph.

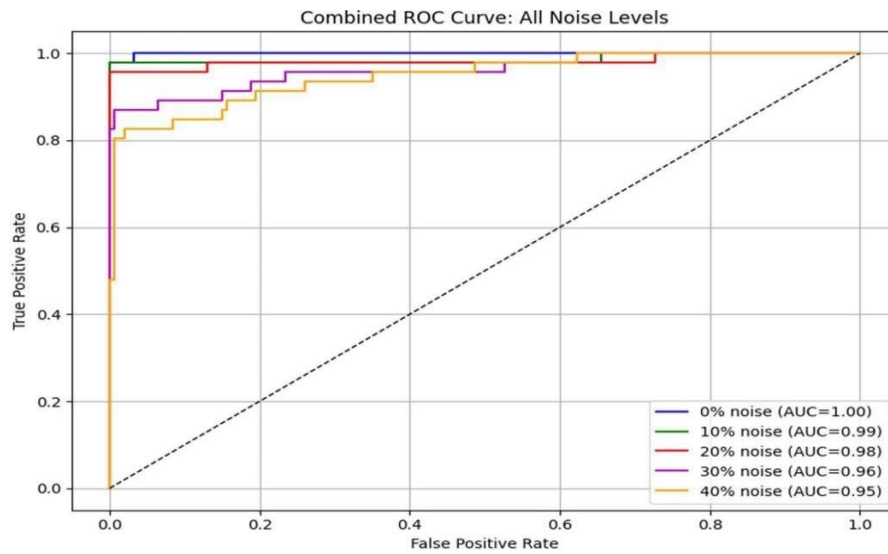


Fig-2: ROC Curves Across Noise Levels for PSO-Optimized Random Forest

(c) Confusion matrix comparison across noise levels

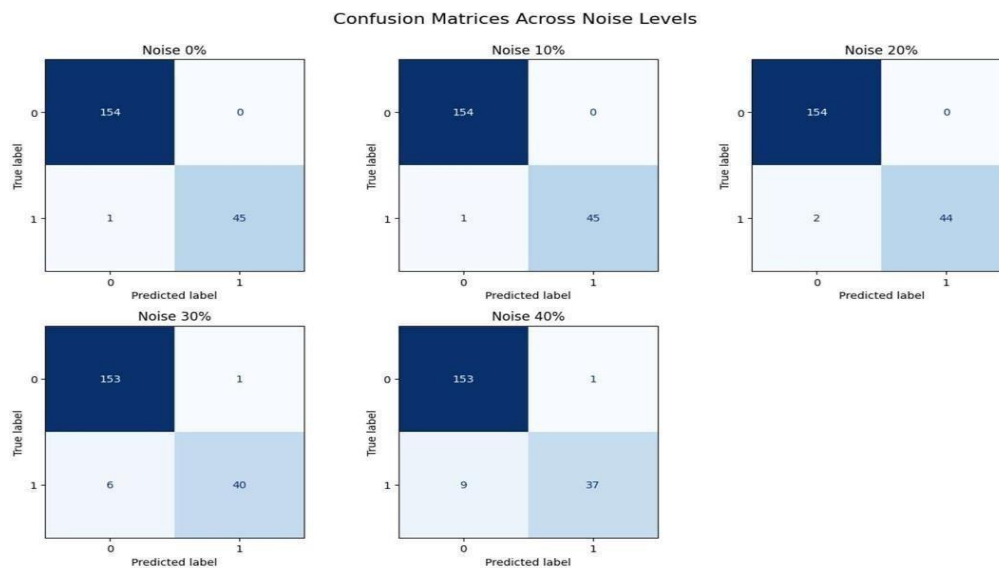


Fig 3. Noise-Level-Wise Confusion Matrix for PSO-Optimized Random Forest

5. Discussion & Conclusion

This study proposed a leakage-aware Particle Swarm Optimization -Random Forest (PSO-RF) framework for prostate cancer survival

prediction using the SEER dataset. The integration of leakage control measures specifically the removal of target-derivative variables ensure that performance estimates reflect true generalization

capability rather than information contamination. Across all experimental settings, the PSO-RF approach demonstrated remarkable resilience to injected Gaussian noise, with Accuracy remaining above 0.95 and F1-score above 0.92 up to 30% noise (Figure 1). Even at the highest noise level (40%), performance degradation was gradual, suggesting that PSO-driven hyperparameter tuning confirms robustness by optimizing model complexity for each fold and noise scenario.

The combined ROC curves (Figure 2) further reinforce this conclusion, with the Area Under the Curve (AUC) staying consistently above 0.95 for noise levels up to 40%. This indicates that the PSO-RF classifier retains strong discriminative capability under substantial feature perturbations.

The confusion matrix grid (Figure 3) provides a detailed view of error patterns across noise intensities. Misclassifications increase slightly as noise grows, but the true positive rate remains dominant even at higher noise levels critical for medical applications where failing to detect at-risk patients could have severe consequences.

From a methodological perspective, two factors underpin these results:

1. Leakage-aware feature filtering, which enhances validity and reproducibility by eliminating spurious correlations.

2. Noise-level specific optimization via PSO, which tunes hyperparameters dynamically to mitigate noise-induced variance.

From a clinical informatics perspective, these findings are highly relevant. Hospital datasets frequently suffer from missing values, coding inconsistencies and measurement noise. A predictive model that retains high performance in the face of such imperfections is more likely to deliver reliable results in real-world deployment. However, limitations remain. The current study treated survival as a binary classification problem rather than performing time-to-event survival analysis. Moreover, while SEER data is comprehensive, it reflects a U.S only patient population, external validation on multi-institutional international datasets would strengthen generalizability.

In conclusion, the proposed leakage-aware PSO-RF pipeline not only addresses the long-standing problem of hidden data leakage in medical AI but also demonstrates high noise robustness across realistic perturbation levels. This dual novelty combining methodological rigor with resilience positions the framework as a strong candidate for deployment in real-world prostate cancer survival prediction systems and as a replicable template for other medical prediction tasks involving complex, imperfect datasets.

References:

1. Aguerchi, K., Jabrane, Y., Habba, M. and et al. 2024. A CNN hyperparameters optimization based on particle swarm optimization for mammography breast cancer classification. *J. Imaging*. 10: 30.
2. Alsaykhan, L.K. and et al. 2024. A hybrid detection model for acute lymphocytic leukemia using support vector machine and particle swarm optimization. *Sci. Rep.* 14: Article 24483.
3. Christodoulou, E., Ma, J., Collins, G.S., Steyerberg, E.W., Verbakel, J.Y. and Van Calster, B. 2019. A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *J. Clin. Epidemiol.* 110: 12–22.
4. Deo, R.C. 2015. Machine learning in medicine. *Circulation*. 132(20): 1920–1930.
5. Kaufman, S., Rosset, S., Perlich, C. and Stitelman, O. 2012. Leakage in data mining: Formulation, detection, and avoidance. *ACM Trans. Knowl. Discov. Data*. 6(4): 15.
6. Kazerani, R. 2024. Improving breast cancer diagnosis accuracy by particle swarm optimization feature selection. *Discov. Appl. Sci.*
7. Molaei, S., Cirillo, S. and Solimando, G. 2024. Cancer detection using a new hybrid method based on pattern recognition in microRNAs combining particle swarm optimization algorithm and artificial neural network.
8. Steyerberg, E.W. and et al. 2013. Prognosis research strategy (PROGRESS)3: Prognostic model research. *PLoS Med.* 10(2): e1001381.
9. Subbaswamy, A. and Saria, S. 2020. From development to deployment: Dataset shift, causality, and shift- stable models in health AI. *Biostatistics*. 21(2): 345–352.
10. Van Gerven, M. and Bohte, S. 2017. Editorial: Artificial neural networks as models of neural information processing. *Front. Comput. Neurosci.* 11: Article 114.
11. Ying, Y., Castro, J., Dodd, J., Singh, K. and Hines, K. 2018. Particle swarm optimization for parameter determination and feature selection of support vector machines. *Expert Syst. Appl.* 36(5): 8898–8908