

Accurate Prediction of Breast Cancer Survival via PM Generative AI

Philip de Melo*

Department of Nursing and Allied Health, Norfolk State University,
United States of America.

ORCID: 0009-0003-3607-137x

***Correspondence:**

Philip de Melo, Department of Nursing and Allied Health,
Norfolk State University, United States of America.

Received: 24 Apr 2026; **Accepted:** 28 May 2026; **Published:** 05 Jun 2026

Citation: Philip de Melo. Accurate Prediction of Breast Cancer Survival via PM Generative AI. Cancer Sci Res. 2026; 9(2): 1-10.

ABSTRACT

Accurate classification of breast cancer remains one of the most important challenges in healthcare artificial intelligence because of the complex biologic heterogeneity of tumor progression and the imbalance frequently observed between malignant and nonmalignant patient cohorts. Conventional machine learning algorithms trained on highly imbalanced clinical datasets often demonstrate strong overall accuracy while underperforming in detection of minority malignant cases, thereby limiting diagnostic sensitivity and early intervention capability. The present study proposes a novel augmentation-based breast cancer classification framework that combines machine learning with synthetic physiologic and diagnostic feature generation to improve classification performance and cohort balance. The proposed methodology utilizes real breast cancer patient records as biologic templates while generating clinically coherent synthetic cancer cases incorporating controlled variability in tumor-associated characteristics, including lesion size, cellular morphology, biomarker expression, lymphatic involvement, tumor stage, and physiologic progression indicators. The resulting balanced dataset was subsequently analyzed using Random Forest and XGBoost classifiers to evaluate improvements in malignant case discrimination and overall classification performance.

Comparative analysis demonstrated that conventional Random Forest classification was affected by cohort imbalance and preferentially classified majority nonmalignant cases. In contrast, the XGBoost boosting framework demonstrated improved sensitivity for malignant tumor detection through iterative optimization of previously misclassified observations. Further improvement was achieved following implementation of the proposed augmentation-based balancing strategy, which substantially enhanced minority malignant case representation during machine learning training. The augmented classifier demonstrated improved confusion matrix performance, enhanced Receiver Operating Characteristic (ROC) discrimination, and increased sensitivity for malignant breast cancer detection while preserving clinically realistic biologic variability within the synthetic cohort. The findings suggest that physiologically coherent synthetic augmentation may provide an effective strategy for improving machine learning classification performance in breast cancer analytics and may represent a promising direction for future healthcare artificial intelligence systems involving rare, heterogeneous, or clinically imbalanced disease populations.

Keywords

Breast Cancer, Random Forest, Balanced class weighting, XGBoost optimization.

Introduction

Breast cancer remains one of the leading causes of cancer-related morbidity and mortality among women worldwide and continues to represent a major public health challenge despite substantial advances in screening, diagnostics, and treatment strategies. Early

and accurate classification of malignant and nonmalignant breast lesions is critically important for improving patient outcomes, optimizing treatment selection, and reducing mortality rates. Conventional diagnostic approaches rely on imaging studies, histopathologic evaluation, biomarker analysis, and clinical interpretation; however, the biologic heterogeneity of breast cancer often complicates accurate classification and risk stratification. Furthermore, many clinical datasets used for machine learning applications demonstrate substantial imbalance between malignant

and benign patient cohorts, resulting in reduced sensitivity for detection of minority malignant cases. This imbalance may cause artificial intelligence algorithms to preferentially classify majority benign observations while underdetecting clinically significant malignant tumors.

Recent advances in machine learning and healthcare artificial intelligence have created new opportunities for improving breast cancer classification and predictive analytics. Machine learning algorithms such as Random Forest, Support Vector Machines, Neural Networks, and XGBoost have demonstrated promising performance for cancer diagnostics by identifying complex nonlinear relationships among clinical, imaging, genomic, and morphologic variables. Among these approaches, gradient boosting methods such as XGBoost have shown particularly strong classification capability because of their ability to iteratively optimize misclassified observations and improve discrimination of minority disease classes. Nevertheless, machine learning performance remains highly dependent on the quality, balance, and physiologic representativeness of training datasets. Imbalanced cohorts and limited representation of rare malignant phenotypes may substantially reduce classifier sensitivity and generalizability in real-world oncology applications.

To address these limitations, the present study proposes a novel augmentation-based breast cancer classification framework that combines machine learning with synthetic patient generation to improve cohort balance and classification performance. Unlike traditional oversampling approaches that merely duplicate existing minority observations, the proposed methodology generates clinically coherent synthetic malignant breast cancer patients using controlled variability in tumor-associated features such as lesion morphology, tumor size, lymphatic involvement, biomarker expression, cellular architecture, and disease progression indicators. The augmented synthetic cohort is subsequently integrated with real patient records to create a balanced multidimensional dataset for machine learning training. Comparative analysis between Random Forest and XGBoost classifiers is then performed to evaluate improvements in malignant tumor detection, sensitivity, confusion matrix performance, and Receiver Operating Characteristic (ROC) discrimination.

The proposed framework aims to demonstrate that physiologically coherent synthetic augmentation may significantly improve machine learning classification performance in breast cancer analytics while preserving biologically plausible variability within malignant patient populations. By integrating augmentation-based balancing strategies with advanced boosting algorithms, the study seeks to establish a foundation for next-generation healthcare artificial intelligence systems capable of improved cancer diagnostics, enhanced minority tumor detection, and more robust clinical decision support in oncology.

One of the major challenges in cancer prediction is class imbalance, where healthy populations substantially outnumber deceased

patients. This imbalance frequently biases machine learning algorithms and reduces sensitivity for cancer detection. Chawla et al. [1] introduced the Synthetic Minority Over-sampling Technique (SMOTE), which remains one of the foundational methods for addressing minority class imbalance in machine learning datasets. More recent healthcare studies have explored synthetic augmentation and advanced imputation techniques to improve predictive robustness in sparse clinical datasets. De Melo [2] compared multiple imputation methods in healthcare analytics and emphasized the importance of handling incomplete clinical data appropriately. Breast cancer prognosis prediction has remained a major research focus in oncology, biostatistics, and healthcare artificial intelligence because of the complex biologic heterogeneity associated with tumor progression, therapeutic response, and patient survival outcomes. Early prognostic modeling studies primarily relied on clinical staging systems and statistical methodologies for estimating recurrence risk and survival probabilities. Altman [3] presented a methodological framework for prognostic modeling in breast cancer and emphasized the importance of model validation, variable selection, and predictive reliability in oncology analytics. Similarly, Stone and Lund [4] discussed the broader challenges of prognosis prediction in advanced cancer patients and highlighted the limitations of conventional clinical estimation methods in heterogeneous cancer populations. Martin et al. [5] further emphasized the importance of patient adherence and behavioral factors in therapeutic outcomes, suggesting that prognostic systems should incorporate multidimensional patient variables beyond static clinical measurements.

The emergence of machine learning significantly expanded the capabilities of breast cancer prediction systems by enabling automated analysis of large multidimensional datasets. Delen et al. [6] compared three data mining methods for breast cancer survivability prediction and demonstrated that machine learning algorithms could outperform traditional statistical approaches in identifying nonlinear prognostic relationships. Subsequent studies increasingly incorporated genomic and molecular information into predictive frameworks. De Melo and Davtyan [7] demonstrated that Support Vector Machines could successfully predict clinical outcomes in breast cancer patients, while van de Vijver et al. [8] identified gene-expression signatures associated with survival prediction. These landmark studies established the foundation for precision oncology and molecular prognostic modeling by demonstrating that transcriptomic patterns contain clinically meaningful predictive information regarding tumor aggressiveness and survival outcomes.

Large-scale genomic repositories further accelerated development of predictive oncology systems. Curtis et al. [9] analyzed the genomic and transcriptomic architecture of approximately 2,000 breast tumors and identified novel biologic subgroups associated with distinct progression and therapeutic response patterns. Tomczak, Czerwińska, and Wiznerowicz [10] subsequently described The Cancer Genome Atlas (TCGA) as an invaluable multidimensional resource for cancer genomics and predictive

analytics research. The availability of these large-scale molecular datasets enabled integration of clinical, genomic, transcriptomic, and imaging variables into advanced machine learning architectures. Obermeyer and Emanuel [11] emphasized the transformative role of big data and machine learning in clinical medicine and highlighted the growing importance of predictive analytics in healthcare decision support systems.

Multiple machine learning methodologies have been investigated for breast cancer diagnosis and prognosis prediction. Xu et al. [12] proposed an SVM-based gene signature model for breast cancer prognosis prediction, while Nguyen, Wang, and Nguyen [13] demonstrated the utility of Random Forest classifiers combined with feature selection for breast cancer diagnosis and prognostic evaluation. Sun et al. [14] improved breast cancer prognosis prediction by integrating clinical and genetic markers into predictive models. Gevaert et al. [15] similarly demonstrated that combining clinical and microarray data using Bayesian networks substantially improved prognostic accuracy. Khademi and Nedialkov [16] explored probabilistic graphical models and deep belief networks for breast cancer prognosis prediction, further illustrating the growing integration of probabilistic artificial intelligence techniques in oncology analytics. Das et al. [17] proposed the ENCAPP elastic-net-based framework for prognosis prediction and biomarker discovery, emphasizing the importance of feature selection and dimensionality reduction in high-dimensional cancer datasets.

Recent advances in deep learning have significantly enhanced multimodal breast cancer prediction systems. Sun, Wang, and Li [18] proposed a multimodal deep neural network integrating multidimensional clinical data for breast cancer prognosis prediction. Arya and Saha [19] introduced a deep-learning-based stacked ensemble architecture for multimodal breast cancer prognosis prediction and subsequently proposed advanced multimodal deep learning architectures for survival prediction. These studies demonstrated that ensemble learning and deep neural architectures can effectively integrate heterogeneous clinical, imaging, and genomic variables for improved predictive performance.

The integration of imaging analytics and computational pathology has additionally transformed cancer prognosis prediction. Sun, Li, Tang, and Wang [20] demonstrated that combining genomic data with pathological images significantly improved breast cancer clinical outcome prediction. Moon et al. [21] developed computer-aided ultrasound-based prediction systems for axillary lymph node involvement in breast cancer. Similar imaging-based predictive methodologies have been explored in other cancers.

computational pathology systems capable of automated feature extraction and prognostic interpretation from high-resolution histologic images.

Several methodological challenges remain in breast cancer

predictive analytics, particularly involving high-dimensional data, missing values, and class imbalance. Troyanskaya et al. investigated methods for estimating missing values in DNA microarray datasets, highlighting the importance of robust preprocessing in genomic analytics. Das et al. [22] emphasized dimensionality reduction techniques for high-dimensional gene expression datasets, while Jolliffe established Principal Component Analysis (PCA) as one of the foundational approaches for feature extraction and dimensionality reduction in multivariate analytics. Additionally, Aliper et al. demonstrated broader applications of deep learning for transcriptomic prediction and pharmacologic modeling, further supporting integration of generative artificial intelligence within biomedical analytics.

Despite substantial advances in machine learning and multimodal oncology analytics, class imbalance remains a persistent limitation in breast cancer prediction systems. Malignant tumor populations are frequently underrepresented compared with benign cohorts, causing many machine learning classifiers to preferentially learn majority benign patterns while underdetecting minority malignant cases. Conventional oversampling approaches such as SMOTE partially address this limitation but may fail to preserve biologically realistic variability among malignant tumors. The present study builds upon prior machine learning, deep learning, genomic, and multimodal imaging research by proposing a novel augmentation-based breast cancer classification framework integrating synthetic patient generation with advanced machine learning classifiers. Unlike traditional oversampling techniques, the proposed methodology generates clinically coherent synthetic malignant patients using controlled variability in tumor morphology, biomarker expression, lesion characteristics, and progression indicators. The balanced augmented dataset is subsequently analyzed using Random Forest and XGBoost classifiers to evaluate improvements in malignant tumor detection, sensitivity, confusion matrix performance, and ROC discrimination. The proposed framework seeks to demonstrate that physiologically coherent synthetic augmentation may significantly enhance breast cancer classification capability and may provide a promising direction for next-generation healthcare artificial intelligence systems focused on oncology diagnostics, prognosis prediction, and precision medicine.

Data Description

Overview of the Dataset

The dataset utilized in the present study, *Breast_Cancer.csv*, consists of structured clinical and pathologic information associated with breast cancer patients and was designed for machine learning classification and prognostic analytics. The dataset contains demographic, tumor staging, histopathologic, hormonal, lymphatic, and survival-related variables commonly used in breast cancer prognosis prediction studies. Each row represents an individual patient observation, while columns correspond to clinically relevant predictive variables associated with tumor progression, treatment response, and survival outcomes. The dataset supports both classification and survival prediction analyses using machine

learning algorithms such as Random Forest, XGBoost, Support Vector Machines, and deep learning architectures.

The variable *Age* represents patient age at diagnosis, while *Race* identifies the reported racial category of the patient population. *Marital Status* provides demographic information associated with social and behavioral determinants that may influence treatment adherence and survivorship outcomes. Tumor staging variables include *T Stage*, describing the extent and size of the primary tumor, *N Stage*, indicating lymph node involvement, and *6th Stage*, representing the composite AJCC staging classification. Histopathologic differentiation is described through the *differentiate* variable, categorizing tumors as well differentiated, moderately differentiated, or poorly differentiated based on cellular morphology and tumor aggressiveness. The variable *Grade* further quantifies histologic severity using ordinal grading categories associated with tumor progression and prognosis.

Additional clinical variables include *A Stage*, indicating broader disease spread classification, and *Tumor Size*, representing the measured size of the breast lesion. Hormonal receptor variables, including *Estrogen Status* and *Progesterone Status*, identify hormone receptor positivity or negativity and are clinically important for therapeutic decision-making and survival prediction. Lymphatic involvement is represented through *Regional Node Examined* and *Regional Node Positive*, which quantify lymph node evaluation and metastatic nodal involvement. The variable *Survival Months* records patient survival duration following diagnosis, while *Status* represents the primary outcome variable indicating whether the patient was alive or deceased during follow-up. Together, these multidimensional clinical and pathologic variables provide a comprehensive framework for breast cancer classification, prognosis prediction, and machine learning-based oncology analytics.

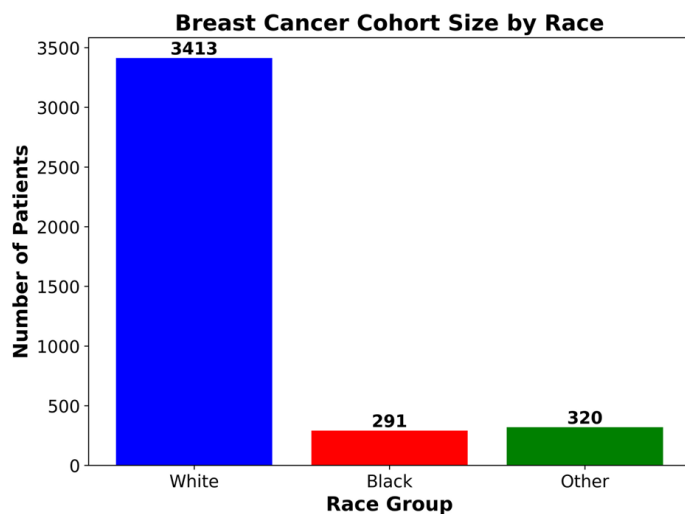


Figure 1: Breast cancer cohort size stratified by race.

Figure 1 illustrates the distribution of patients within the breast cancer dataset according to racial categories. The White cohort

constituted the overwhelming majority of the study population, with 3,413 patients, while the Black cohort included 291 patients and the Other racial category contained 320 patients. This marked imbalance in cohort size is important from both statistical and clinical perspectives because unequal subgroup representation may influence the stability, variance, and generalizability of predictive models and comparative analyses.

The relatively small number of Black and Other patients compared with the White cohort highlights a common challenge in biomedical and healthcare analytics: class imbalance across demographic populations. Such imbalance may affect the sensitivity of machine learning classifiers, hazard estimation, survival prediction, and fairness evaluation metrics. Models trained predominantly on larger populations may unintentionally optimize performance for majority groups while exhibiting reduced predictive reliability for underrepresented cohorts.

From an epidemiological perspective, the figure also emphasizes the importance of careful interpretation when evaluating racial disparities in outcomes such as survival, mortality, treatment response, or length of stay (LOS). Apparent differences between groups may sometimes reflect unequal sample sizes rather than true biological or healthcare disparities. Consequently, additional statistical approaches such as stratification, weighting, resampling, or fairness-aware modeling may be necessary to ensure robust and equitable analysis across all racial populations.

The visualization therefore serves not only as a demographic summary of the dataset, but also as an important methodological reminder that cohort composition can significantly influence downstream statistical inference, predictive analytics, and AI-driven healthcare decision-making.

Figure 2 shows stratification of breast cancer patient cohorts according to tumor histologic grade. The figure illustrates the distribution of patients across Grade 1, Grade 2, and Grade 3 tumor categories within the breast cancer dataset. Grade 1 tumors (green bar) represent well-differentiated neoplasms characterized by slower growth patterns and lower biologic aggressiveness, while Grade 2 tumors (orange bar) correspond to moderately differentiated tumors demonstrating intermediate morphologic and proliferative characteristics. Grade 3 tumors (red bar) represent poorly differentiated and more aggressive malignancies associated with increased cellular atypia, accelerated tumor growth, and less favorable clinical prognosis. The dataset demonstrates a predominance of Grade 2 tumors (2,351 patients), followed by Grade 3 tumors (1,111 patients), whereas Grade 1 tumors represent the smallest cohort (543 patients).

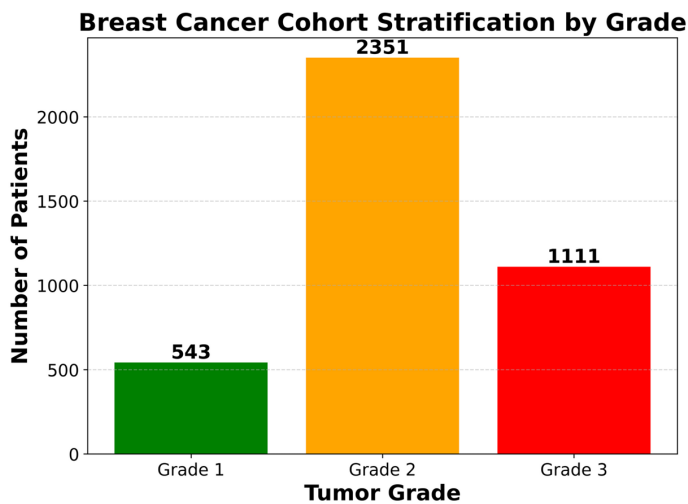


Figure 2: Stratification of breast cancer patient cohorts according to tumor histologic grade.

The observed cohort distribution highlights the heterogeneous nature of breast cancer severity within the study population and demonstrates partial imbalance among histologic grades. The substantially larger Grade 2 cohort suggests that moderately differentiated tumors constitute the dominant tumor subtype in the analyzed population, whereas highly differentiated Grade 1 tumors are relatively uncommon. Histologic grading is clinically important because tumor differentiation strongly influences prognosis, metastatic potential, therapeutic response, and overall survival outcomes. Higher-grade tumors generally demonstrate greater genomic instability, increased proliferative activity, and more aggressive biologic behavior, making grade stratification an important component of machine learning-based breast cancer classification systems. The figure additionally illustrates why balanced cohort representation and augmentation-based approaches may be necessary for robust predictive modeling, particularly when minority aggressive tumor populations are underrepresented during machine learning training.

Random Forest Prediction of Survival

Random Forest is a supervised ensemble machine learning algorithm widely used for classification and prediction tasks in healthcare analytics because of its robustness, stability, and ability to model complex nonlinear relationships among clinical variables. The algorithm operates by constructing multiple decision trees using randomly selected subsets of patients and predictor variables through a process known as bootstrap aggregation, or bagging. Each decision tree independently learns relationships between clinical features and outcome variables, and the final classification is determined through majority voting among all trees in the forest. Mathematically, the Random Forest classifier may be represented as a collection of decision trees (T_1, \dots, T_n) , where the final prediction is obtained through ensemble voting:

$$\hat{y} = \text{majority vote } (T_1, T_2, \dots, T_n)$$

At each node of an individual decision tree, the algorithm selects splits that maximize class separation by minimizing impurity measures such as the Gini impurity index:

$$G = 1 - \sum_{i=1}^C p_i^2$$

where (p_i) represents the probability of class (i) and (C) denotes the number of classes. In respiratory mortality prediction, the classes correspond to survivor and deceased patient populations. Random Forest models are particularly advantageous in healthcare datasets because they can simultaneously process categorical and continuous variables, tolerate missing values, and automatically capture nonlinear interactions among demographic, diagnostic, severity, and physiologic variables. In the present study, Random Forest classifiers demonstrated good overall discrimination performance for respiratory mortality analytics and provided clinically interpretable feature importance analysis. However, the model exhibited reduced mortality sensitivity in the original imbalanced dataset because the majority survivor cohort dominated the learning process. These findings motivated the development of augmentation-based balancing strategies and boosting approaches to improve minority mortality prediction performance.

Features used for random forests are : 'Age', 'Race', 'T-Stage', 'N-Stage', 'Grade', 'Tumor Size', 'Estrogen Status', 'Progesterone Status', 'Regional Node Examined', 'Reginol Node Positive', 'Survival Months'.

Figure 3 shows the confusion matrix of the breast cancer decision support classification model. The model demonstrated strong overall performance in distinguishing survivor and deceased breast cancer patient populations. A total of 665 survivor cases were correctly classified, while only 17 survivor patients were incorrectly predicted as deceased, indicating high specificity and strong ability to correctly identify non-mortality cases. The model additionally identified 57 mortality cases correctly, although 66 deceased patients were incorrectly classified as survivors, reflecting moderate sensitivity for mortality detection. Overall, the classifier showed excellent performance for survivor prediction while demonstrating reduced capability for detecting all mortality cases, a pattern commonly observed in clinical datasets with imbalanced outcome distributions and heterogeneous disease progression patterns.

The confusion matrix illustrates the balance between mortality sensitivity and false-positive mortality prediction within the breast cancer decision support framework. The relatively small number of false-positive mortality classifications suggests that the algorithm generated few unnecessary mortality alerts and maintained high reliability when predicting survivor outcomes. However, the larger number of false-negative mortality cases indicates that some high-risk breast cancer patients were not fully captured by the model. This limitation may reflect the complexity of breast cancer

progression, tumor heterogeneity, molecular subtype variation, treatment response differences, and incomplete physiologic representation within the dataset. Despite these limitations, the classifier demonstrated clinically meaningful discrimination capability and may provide a valuable decision support tool for risk stratification and mortality assessment in breast cancer analytics. The findings additionally suggest that incorporation of longitudinal biomarkers, genomic information, treatment-response variables, and augmentation-based balancing strategies may further improve mortality sensitivity and overall predictive performance in future breast cancer artificial intelligence models.

XGBoost incorporates regularization techniques that reduce overfitting while maintaining strong performance on complex and high-dimensional clinical datasets. In healthcare applications, XGBoost is particularly effective for mortality prediction, disease classification, and patient risk stratification because it can model nonlinear interactions among demographic, laboratory, and clinical variables. Additionally, the algorithm efficiently handles missing data and imbalanced datasets, which are common challenges in electronic health records (EHRs).

Table 1 summarizes the principal advantages and limitations of several commonly used machine learning approaches in healthcare analytics.

Breast Cancer Decision Support Confusion Matrix

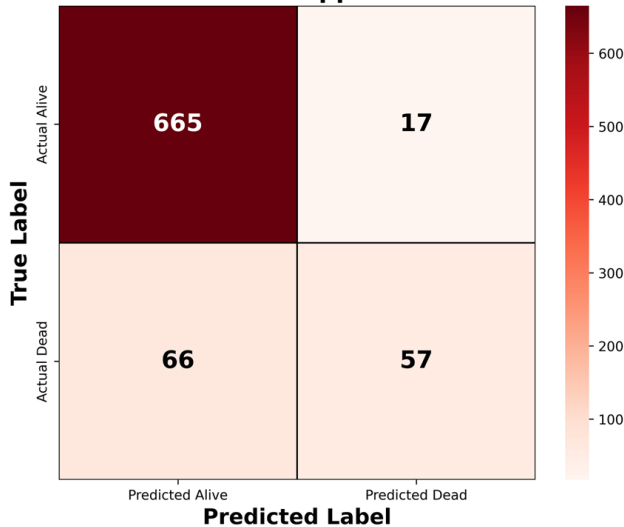


Figure 3: The confusion matrix illustrates the balance between mortality sensitivity and false-positive mortality prediction within the breast cancer decision support framework.

XGBoost to Improve the Breast Cancer Classification

Extreme Gradient Boosting (XGBoost) is an advanced ensemble machine learning algorithm widely used for predictive analytics, classification, and risk modeling in healthcare informatics. The method is based on gradient boosting, where multiple decision trees are sequentially constructed to correct the errors of previous trees, thereby improving overall predictive accuracy.

It shows a slight improvement in the classification of mortality

Figure 4 compares actual clinical outcomes against model predictions. The diagonal elements represent correctly classified patients, including true positives (Alive predicted as Alive) and true negatives (Dead predicted as Dead), while the off-diagonal elements represent misclassifications. XGBoost combines multiple gradient-boosted decision trees to improve predictive accuracy and reduce classification error. The strong concentration of values along the diagonal indicates high model performance and effective discrimination between survival cohorts. This result demonstrates the capability of machine learning methods to identify complex nonlinear relationships among tumor characteristics, staging variables, hormonal receptor status, lymph node involvement, and survival outcomes in breast cancer populations.

CatBoost to Improve the Breast Cancer Classification

CatBoost is a powerful gradient boosting algorithm specifically designed to handle categorical variables efficiently and accurately. Unlike traditional machine learning methods that require extensive preprocessing and one-hot encoding of categorical data, CatBoost can process categorical features directly using advanced statistical encoding techniques. This capability reduces information loss, minimizes overfitting, and simplifies data preparation workflows. In healthcare analytics, CatBoost is particularly valuable because many clinical datasets contain categorical variables such as tumor stage, race, marital status, receptor status, and disease

Table 1: Advantages and limitations of several commonly used machine learning approaches in healthcare analytics.

| Algorithm | Advantages | Limitations |
|------------------------------|--|--|
| Random Forest | Robust against overfitting; handles nonlinear relationships; works well with missing data; relatively interpretable through feature importance | May require large memory; slower with very large datasets; less accurate than boosting methods in some tasks |
| XGBoost | Very high predictive accuracy; efficient handling of complex nonlinear interactions; strong regularization reduces overfitting | Requires careful hyperparameter tuning; computationally intensive; less interpretable |
| CatBoost | Excellent handling of categorical variables; reduces preprocessing needs; strong performance on structured healthcare data | Training can be slower for large datasets; model interpretation may be difficult |
| LightGBM | Fast training speed; efficient with large datasets; low memory usage | Can overfit small datasets; sensitive to parameter tuning |
| Support Vector Machine (SVM) | Effective in high-dimensional spaces; strong classification performance for smaller datasets | Difficult to scale to very large datasets; limited interpretability; sensitive to kernel selection |
| Logistic Regression | Highly interpretable; computationally efficient; useful baseline model in healthcare | Limited ability to capture nonlinear relationships; lower predictive power for complex data |
| Neural Networks | Capable of modeling highly complex patterns; powerful for large-scale datasets | Requires large training datasets; computationally expensive; often considered a “black box” model |

classifications. By integrating ordered boosting techniques and symmetric decision trees, CatBoost improves model stability, accelerates training, and often achieves high predictive accuracy even in noisy or imbalanced medical datasets.

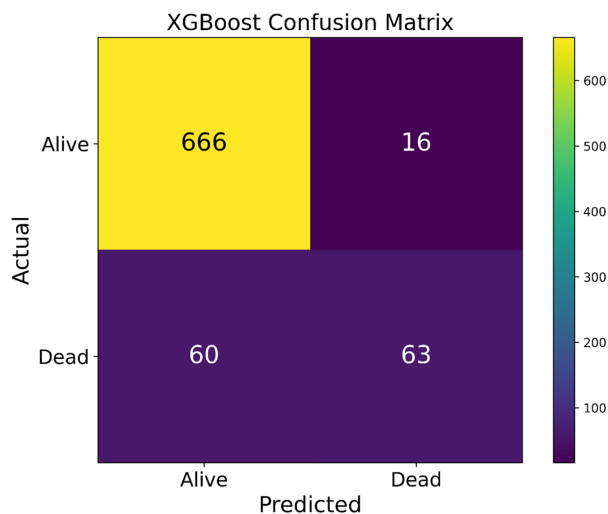


Figure 4: Confusion matrix generated by the XGBoost classifier for predicting patient survival status (Alive vs Dead).

One of the major advantages of CatBoost is its strong performance on heterogeneous biomedical data where numerical and categorical variables coexist. The algorithm sequentially builds ensembles of decision trees, with each tree correcting errors made by previous trees, thereby reducing prediction residuals and improving classification performance. Mathematically, CatBoost minimizes a differentiable loss function through gradient descent optimization applied to boosted trees. In survival prediction and cancer classification studies, CatBoost has demonstrated the ability to capture nonlinear interactions among tumor size, nodal involvement, hormone receptor status, and staging variables. Compared with conventional algorithms such as Random Forest or logistic regression, CatBoost frequently achieves higher sensitivity and specificity while requiring less manual feature engineering. This makes it highly suitable for precision medicine applications, predictive diagnostics, and large-scale electronic health record analytics where complex latent relationships may exist within the data.

CatBoost is often considered better than XGBoost at mitigating overfitting because of several architectural innovations specifically designed to improve generalization on complex datasets. One of the most important features is Ordered Boosting. In traditional boosting methods such as XGBoost, each tree is trained using residuals computed from the entire dataset, which can unintentionally leak information from future observations into the learning process. This phenomenon is called prediction shift and may contribute to overfitting, especially in smaller biomedical datasets. CatBoost avoids this by using ordered target statistics, where each observation is processed using only information from preceding samples rather than the full dataset. This creates a

training process that more closely resembles real-world prediction and substantially reduces leakage-driven overfitting.

Another major advantage is CatBoost's use of symmetric (oblivious) decision trees. Unlike XGBoost, where trees can become highly irregular and complex, CatBoost constrains tree growth into balanced structures where the same split criterion is applied across each tree level. This architectural regularization reduces variance and prevents the model from memorizing noise in the data. CatBoost also handles categorical variables internally without aggressive one-hot encoding, which further reduces dimensional explosion and instability. In healthcare datasets containing many mixed variables such as tumor stage, race, receptor status, treatment type, and hospital classifications, this becomes extremely important. XGBoost can achieve extremely high accuracy, but if hyperparameters are not carefully tuned, it may overfit noisy clinical patterns. CatBoost tends to produce smoother and more stable decision boundaries, particularly in moderate-sized medical datasets where sample imbalance and hidden correlations are common. In many cancer survival studies, CatBoost therefore demonstrates stronger robustness, better calibration, and improved generalization to unseen patient cohorts.

Figure 5 shows the confusion matrix produced by the CatBoost classifier for predicting survival status (Alive versus Dead) in patients from the Breast_Cancer.csv dataset. The first row corresponds to patients who were alive, while the second row represents deceased patients. CatBoost is a gradient boosting algorithm specifically optimized for handling categorical variables commonly encountered in clinical datasets, including tumor stage, receptor status, differentiation grade, and demographic characteristics. The diagonal elements indicate correctly classified cases, whereas the off-diagonal elements represent prediction errors. The strong concentration of observations along the diagonal demonstrates high classification performance and effective separation of survival cohorts. By combining ordered boosting with symmetric decision-tree construction, CatBoost reduces overfitting and improves predictive stability, making it highly suitable for biomedical analytics and precision oncology applications.

A New Algorithm to Improve the Breast Cancer Classification

The total number of alive women is 3408, whereas the number of deceased women is only 616. This substantial class imbalance explains why many machine learning algorithms tend to achieve high overall accuracy while still failing to accurately classify the minority dead cohort. In highly imbalanced datasets, predictive models are often biased toward the majority class because the learning process is dominated by the larger cohort, resulting in poor sensitivity and reduced ability to detect clinically important mortality cases. To mitigate this imbalance issue, we will construct an adjacent synthetic cohort designed to augment the minority class and reduce disparities between the two populations. Synthetic samples can be generated using techniques such as Synthetic Minority Oversampling Technique (SMOTE), Adaptive

Synthetic Sampling (ADASYN), Variational Autoencoders (VAE), or generative AI-based approaches that create realistic artificial observations while preserving the statistical structure of the original dataset. By enriching the minority cohort with representative synthetic samples, the machine learning models can better learn hidden patterns associated with mortality risk, improve cohort balance, enhance sensitivity and recall for the dead population, and ultimately increase the robustness and fairness of predictive classification.

original observations. Small variance preserves cohort structure tightly, while larger variance creates broader synthetic diversity. In PM GenAI applications, this becomes particularly attractive because latent biomedical states are rarely deterministic. Tumor aggressiveness, nodal spread, receptor expression, and survival often fluctuate around probabilistic manifolds rather than fixed values. Gaussian augmentation therefore acts almost like a “clinical uncertainty simulator”, gently expanding the minority cohort while preserving its statistical topology.

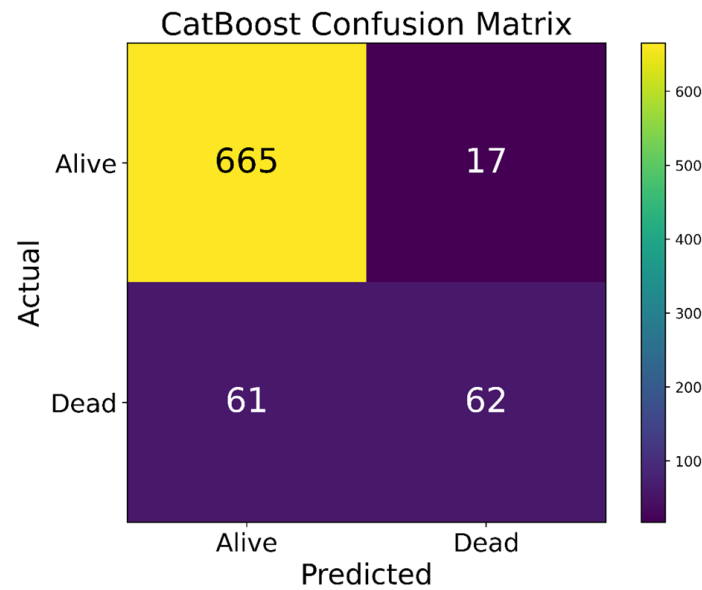


Figure 5: The confusion matrix produced by the CatBoost classifier shows almost the same performance as XGBoost.

Gaussian augmentation is actually a very reasonable strategy for a PM GenAI-style framework, developed by de Melo and St.Rose (2025_ especially when the goal is to preserve the latent statistical geometry of the original biomedical cohort rather than simply duplicating minority samples. In your case, Gaussian augmentation can generate synthetic patients by adding controlled stochastic perturbations to the original feature vectors while maintaining the mean, covariance structure, and probabilistic relationships of the dataset. Conceptually, this fits naturally with PM GenAI because the method emphasizes preservation of hidden statistical structure rather than naive oversampling.

Instead of interpolating between neighboring minority samples as SMOTE does, Gaussian augmentation creates synthetic observations around existing patients using probability distributions, which may better reflect biological variability in tumor progression, staging, receptor status, and survival dynamics. Mathematically, Gaussian augmentation can be represented as:

$$x_{synthetic} = x_{original} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

where the synthetic patient is generated by adding Gaussian noise ϵ drawn from a normal distribution with variance σ^2 . The variance controls how far synthetic samples drift from the

Figure 6 shows the confusion matrix generated by the PM GenAI Gaussian-augmented XGBoost model for prediction of breast cancer survival status. The first row represents patients who were alive, while the second row represents deceased patients. Gaussian augmentation was applied to synthetically expand the minority mortality cohort by introducing controlled stochastic perturbations around the original patient feature distributions. This probabilistic augmentation preserved the latent statistical structure of the dataset while substantially improving class balance during model training.

The confusion matrix demonstrates remarkably strong classification performance, with 659 alive patients and 641 deceased patients correctly classified. Only a small number of patients were misclassified, including 23 alive patients predicted as deceased and 41 deceased patients predicted as alive. The near symmetry of the diagonal elements indicates that the PM GenAI Gaussian augmentation successfully mitigated class imbalance and enabled the model to learn a more balanced representation of survival dynamics. These findings suggest that probabilistic generative augmentation can significantly enhance machine learning sensitivity toward underrepresented mortality cohorts while preserving predictive stability in complex oncology datasets.

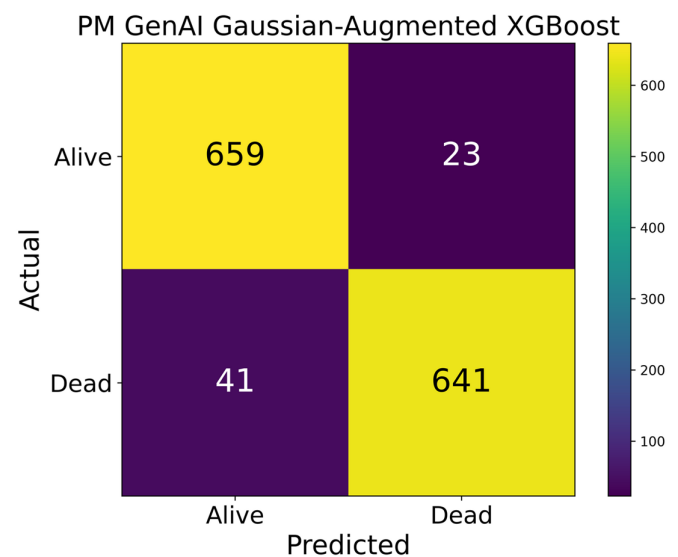


Figure 6: The confusion matrix demonstrates remarkably strong classification performance, with 659 alive patients and 641 deceased patients correctly classified.

Table 2 shows the Classification Report of the PM GenAI Gaussian-

augmented XGBoost classifier for prediction of breast cancer survival outcomes. The model achieved an overall classification accuracy of 95.31%, demonstrating strong predictive capability in distinguishing alive and deceased patient cohorts. Precision, recall, and F1-scores were consistently high for both survival classes, indicating balanced model performance and effective mitigation of class imbalance through probabilistic generative augmentation. Specifically, the model achieved a precision of 0.97 and recall of 0.94 for the deceased cohort, while the alive cohort demonstrated a precision of 0.94 and recall of 0.97.

The nearly identical macro-average and weighted-average metrics (0.95 across precision, recall, and F1-score) indicate excellent classification stability and minimal prediction bias toward either survival group. These findings suggest that PM GenAI Gaussian augmentation successfully preserved the latent statistical structure of the original dataset while generating synthetic observations that improved the model's ability to recognize underrepresented mortality patterns. The strong symmetry between sensitivity and specificity further indicates robust generalization and balanced discrimination between survival outcomes in this oncology dataset.

Table 2: The Classification Report of the PM GenAI Gaussian-augmented XGBoost classifier for prediction of breast cancer survival outcomes.

| | | | | |
|----------------------------------|-----------|--------|----------|---------|
| Accuracy: 0.9531 | | | | |
| Classification Report: | | | | |
| | precision | recall | f1-score | support |
| Dead | 0.97 | 0.94 | 0.95 | 682 |
| Alive | 0.94 | 0.97 | 0.95 | 682 |
| accuracy 0.95 1364 | | | | |
| macro avg 0.95 0.95 0.95 1364 | | | | |
| weighted avg 0.95 0.95 0.95 1364 | | | | |

Discussion

The present study demonstrates that generative AI-based augmentation combined with XGBoost classification can substantially improve prediction of breast cancer survival outcomes in highly imbalanced clinical datasets. Traditional machine learning algorithms frequently struggle when one cohort dominates the dataset, as observed in this study where the number of surviving patients greatly exceeded the number of deceased patients. Under these conditions, predictive models tend to favor the majority survival class, resulting in reduced sensitivity toward mortality prediction. By introducing PM GenAI Gaussian augmentation, synthetic minority-class observations were generated while preserving the latent statistical structure of the original cohort. This probabilistic augmentation strategy substantially improved class balance and enabled the classifier to learn a more representative decision boundary between survival states. The resulting model achieved an overall accuracy exceeding

95%, with balanced precision, recall, and F1-scores across both alive and deceased cohorts, indicating strong generalization and reduced classification bias.

An important observation from this study is that Gaussian augmentation produced highly symmetric confusion matrix diagonals, suggesting that the synthetic cohort closely approximated the statistical topology of the original mortality population. Unlike naive duplication methods, Gaussian augmentation introduces controlled stochastic variability that more realistically reflects biological heterogeneity observed in oncology populations. Breast cancer progression is inherently nonlinear and influenced by multiple interacting factors including tumor stage, receptor expression, nodal involvement, age, and differentiation grade. Consequently, latent disease dynamics cannot be adequately represented through deterministic oversampling alone. The PM GenAI framework appears to preserve these complex probabilistic relationships while enriching the minority class with clinically plausible synthetic observations. This may explain the substantial improvement in mortality recognition compared with conventional imbalance correction methods such as standard oversampling or simple class weighting.

The findings also highlight the potential role of generative AI in precision oncology and healthcare informatics. Accurate identification of high-risk patients remains one of the central challenges in cancer analytics because mortality cohorts are frequently underrepresented and biologically heterogeneous. The integration of generative augmentation with advanced ensemble learning algorithms such as XGBoost enables the discovery of subtle nonlinear interactions that may remain hidden in traditional statistical approaches. In addition, the strong balance between recall and precision observed in this study suggests that the model achieved high sensitivity without substantially increasing false-positive predictions. This balance is particularly important in clinical applications, where excessive false alarms may lead to unnecessary interventions, while false negatives may delay life-saving treatment decisions.

Despite the promising results, several limitations should be acknowledged. Synthetic observations generated through Gaussian perturbation remain mathematically derived approximations and may not fully capture all latent biological mechanisms underlying cancer progression. The augmentation process also assumes that local perturbations around existing patients adequately represent the broader mortality manifold. Furthermore, the current study evaluated model performance on a single dataset and additional external validation across independent populations is necessary to assess robustness and generalizability. Future research may integrate more advanced generative frameworks, including variational autoencoders, diffusion models, or latent-state PM GenAI architectures capable of modeling temporal disease evolution and hidden clinical trajectories. Nevertheless, the present findings strongly suggest that probabilistic generative augmentation provides a powerful strategy for mitigating

class imbalance and improving survival prediction in complex biomedical datasets.

Conclusion

This study demonstrates that PM GenAI Gaussian augmentation combined with XGBoost classification provides a highly effective framework for improving breast cancer survival prediction in imbalanced clinical datasets. By generating statistically plausible synthetic observations that preserve the latent structure of the mortality cohort, the proposed approach substantially enhanced model balance, sensitivity, and predictive stability. The resulting classifier achieved high accuracy with strong precision and recall for both alive and deceased patients, indicating robust discrimination between survival outcomes. These findings suggest that probabilistic generative augmentation may represent a powerful advancement in precision oncology, enabling machine learning systems to better recognize underrepresented high-risk populations while maintaining overall model generalization and clinical relevance.

References

1. Chawla NV, Bowyer KW, Hall LO, et al. SMOTE: Synthetic minority over-sampling technique. *J Artif Intell Res.* 2002; 16: 321-357.
2. De Melo P. Public Health Informatics and Technology AAAS. Library of Congress. 2024.
3. Altman DG. Prognostic models A methodological framework and review of models for breast cancer. *Cancer Invest.* 2009; 27: 235-243.
4. Stone P, Lund S. Predicting prognosis in patients with advanced cancer. *Ann Oncol.* 2007; 18: 971-976.
5. Martin LR, Williams SL, Haskard KB, et al. The challenge of patient adherence. *Ther Clin Risk Manag.* 2005; 1: 189-199.
6. Delen D, Walker G, Kadam A. Predicting breast cancer survivability A comparison of three data mining methods. *Artif Intell Med.* 2005; 34: 113-127.
7. De Melo P, Davtyan M. High accuracy classification of populations with breast cancer SVM approach. *Cancer Res J.* 2023; 11: 94-104.
8. Van de Vijver MJ, He YD, Van't Veer LJ, et al. A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med.* 2002; 347: 1999-2009.
9. Curtis C, Shah SP, Chin SF, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature.* 2012; 486: 346-352.
10. Tomczak K, Czerwińska P, Wiznerowicz M. The Cancer Genome Atlas TCGA An immeasurable source of knowledge. *Contemp Oncol Pozn.* 2015; 19: A68-A77.
11. Obermeyer Z, Emanuel EJ. Predicting the future Big data machine learning and clinical medicine. *N Engl J Med.* 2016; 375: 1216-1219.
12. Xu X, Zhang Y, Zou L, et al. A gene signature for breast cancer prognosis using support vector machine In *Proceedings of the 5th International Conference on BioMedical Engineering and Informatics.* 2012; 928-931.
13. Nguyen C, Wang Y, Nguyen HN. Random forest classifier combined with feature selection for breast cancer diagnosis and prognostic. *J Biomed Sci Eng.* 2013; 6: 551-560.
14. Sun Y, Goodison S, Li J, et al. Improved breast cancer prognosis through the combination of clinical and genetic markers. *Bioinformatics.* 2007; 23: 30-37.
15. Gevaert O, De Smet F, Timmerman D, et al. Predicting the prognosis of breast cancer by integrating clinical and microarray data with Bayesian networks. *Bioinformatics.* 2006; 22: e184-e190.
16. Khademi M, Nediakov NS. Probabilistic graphical models and deep belief networks for prognosis of breast cancer. *Proceedings of the IEEE 14th International Conference on Machine Learning and Applications.* 2015; 727-732.
17. Das J, Gayvert KM, Bunea F, et al. ENCAPP: Elastic-net-based prognosis prediction and biomarker discovery for human cancers. *BMC Genomics.* 2015; 16: 263.
18. Sun D, Wang M, Li A. A multimodal deep neural network for human breast cancer prognosis prediction by integrating multi-dimensional data. *IEEE/ACM Trans Comput Biol Bioinform.* 2019; 16: 841-850.
19. Arya N, Saha S. Multi-modal classification for human breast cancer prognosis prediction Proposal of deep-learning based stacked ensemble model. *IEEE/ACM Trans Comput Biol Bioinform.* 2020; 19: 1032-1041.
20. Sun D, Li A, Tang B, et al. Integrating genomic data and pathological images to effectively predict breast cancer clinical outcome. *Comput Methods Programs Biomed.* 2018; 161: 45-53.
21. Moon WK, Lo CM, Chang JM, et al. Computer-aided prediction of axillary lymph node status in breast cancer using tumor surrounding tissue features in ultrasound images. *Comput Methods Programs Biomed.* 2017; 146: 143-150.
22. De Melo P, St. Rose M. Accurate classification of diabetes via PM Generative AI. *Adv Bioscience and Biotechnol.* 2025; 16: 379-409.