# Chemotherapy Resistance Prediction for Breast Cancer Using Synthetic Multi-Omics Data

Mohamed Ayoub Khabiry
*ITCS*
*Nile University*
Giza, Egypt
M.ayoub2228@nu.edu.eg

## ABSTRACT

Chemotherapy resistance is a major problem in breast cancer treatment that leads to poor survival and quality of life for patients. Here, we leverage synthetic multi-omics data (genomic, transcriptomic, and proteomic features) to predict resistance to chemotherapy using machine learningmodels. Given different models that were tested including Logistic Regression, Support Vector Machines (SVM), and Multi-Layer Perceptrons (MLP), it was theMLP that performed the best at 98.51% accuracy and perfect 1.00 AUC-ROC score. These results highlight the abilityof machine learning to transform cancer research, demonstrating that multi-omics information can be used to make accurate, actionable predictions. These advances offer the ability forclinicians to more precisely tailor therapies, to spare patients the ravages of medications the patient may not need and their attendant toxicities, and to usher in a new, kinder and gentler, custom-fitted age of cancer care **Index Terms—Machine Learning, Multi-Omics Data, Chemotherapy Resistance, Personalized Medicine, Breast Cancer.**

## I. INTRODUCTION

Among the most prevalent and life-threatening diseases worldwide, breast cancer is a major challenge for healthcare systems, affecting millions of lives annually [1, 2]. Despite advances in therapeutic strategies, chemotherapy remains the mainstay of treatment for breast cancer patients. However, a major obstacle limiting the effectiveness of chemotherapy is a phenomenon known as chemotherapy resistance. This resis- tance diminishes the effectiveness of treatment and impairs patient outcomes, while the disease often progresses further [3, 4]. Understanding the underlying mechanisms of chemotherapy resistance and developing methods to predict its onset are fundamental for personalized medicine and therapy [5]. Chemoresistance originates from genetic, transcriptomic, and proteomic interactions that determine cell responses to therapeutic agents [6, 7]. Most predictive models for chemotherapy resistance are based on traditional clinical and demographic data, which lack the sophistication required to represent the complexity of molecular mechanisms in play [8, 9]. The emergence of novel omic technologies in recent times—of particular relevance to geno-, transcripto-, and proteomic datasets—has presented unique opportunities for deciphering the complex mechanisms of resistance [10]. Multi-omics data offer a panoramic view of the biological landscape with biomarkers and pathways related to resistance identifiable within that system. However, such datasets necessitate sophisticated computational techniques that can handle high-dimensional, heterogeneous data [11]. Machine learning has become a powerful tool in oncology. It allows for the processing of complicated data and identification of patterns that are invisible to classic statistical methods [12]. Applying machine learning algorithms to multi-omics data will lead to predictive models identifying chemotherapy resistance with high precision, and shed light on the underlying molecular mechanisms [13]. This is the first approach using synthetic multi-omics data in order to investigate chemotherapy resistance in breast cancer. The dataset combines genomic, transcriptional, and proteomic fea- tures underlying chemotherapy responses. Synthetic data en- ables controlled experimentation without ethical and practical limitations associated with real clinical data [14]. It provides a simulated environment for testing the efficacy of machine learning algorithms under conditions closely approximating real-life clinical scenarios.By emulating such complex biologi-cal relationships and maintaining balanced class representation, synthetic datasets represent ideal testbeds for training and evaluating predictive algorithms [14]. Herein, machine learning prediction of chemotherapy resistance is demonstrated through the implementation of three models: Logistic Regression, Support Vector Machines (SVM), and Multi-Layer Perceptron (MLP). First, Logistic Regression serves as a baseline owing to its simplicity and interpretability, while SVM is suited to model nonlinear rela- in high- Because of the multi-layer structure of MLP, as well as other deep learning models, capturing multi-dimensional data can be possible since they capture complicated patterns and interactions from inside a dataset [12]. The main objective of this paper is the investigation of predictive performance and efficiency of these models in chemotherapy resistance. Besides, studying the performance of these models would outline strengths and weaknesses [13]. The results are expected to emphasize the utility of machine The current work extends the emergent literature using synthetic data to create and test models in machine learning for translation of multi-omics data into clinical decisions that have broader implications for oncology and personalized medicine. These predictive models would revolutionize cancer care, since clinicians could personalize treatment strategies for each patient to minimize therapeutic

failure and the resulting burden of ineffective treatment [4]. Integration of multi-omics data into These predictive models further enable insight into the biology of cancers and the identification of new biomarkers and therapeutic targets 5. Synthetic data is also a major strength in this work, but synthetic data suffers from a number of disadvantages: it cannot capture real-world variability or confounding variables as well as actual clinical data can 14. However, this work provided a good basis for the development and refinement of predictive models with some framework for clinical translation in the future.

## II. METHODOLOGY

The architecture of the developed predictive framework relies on a synthetic multi-omics dataset in order to model chemotherapy resistance in breast cancer. The genomic, transcriptomic, and proteomic features are combined in the dataset to represent complex biological interactions underlying the treatment response. This architecture integrates data preprocessing, model training and performance evaluation, tion, where the main aim is to build robust and interpretable Models capable of predicting resistance to chemotherapy. The synthetic dataset was preprocessed in depth to prepare it for modeling. The missing values of the geno- Iterative approach was used for handling mic and transcriptomic features. imputation approach, which estimates based on the pattern of the data these missing values. In terms of the categorical variables, simple mode imputation was implemented without losing the main class characteristics. Features of a numerical nature, after imputation, were standardized using a standard scaler to assure zero mean and unit variance. A Min-Max scaler was used to normalize selected features, bringing them into a consistent Fig. 1: Methodology workflow overview. This figure depicts the sources of data, data preprocessing, training of models, and evaluation of this study. range to enhance model performance. These steps were crucial to make the data stable and accurate for downstream modeling. The modeling was planned to assess three machine learning architectures comprising Logistic Regression, Support Vector Machine, and Multi-Layer Perceptron. Logistic Regression was taken as a baseline model and gave a straightforward yet effective structure for binary classification.The SVM was done with the radial basis function, kept the regularization parameter C = 1 in order to model nonlinear relationships between the features. Finally, the MLP was a deep learning model that consisted of two hidden layers of 16 and 8 neurons, respectively, using the ReLU activation function. The MLP was trained by using the The learning rate of the Adam optimizer was dynamically adjusted for optimum convergence. Early stopping with a fraction of 20% validation was done to prevent overfitting so that it could generalize well on unseen data.

Interpretability and explainability were also considered. Although deep learning models, such as MLP, perform very well in terms of prediction, their "black-box" nature is one of the reasons for their low clinical acceptance. Therefore, insights gained from logistic regression and SVM, which are
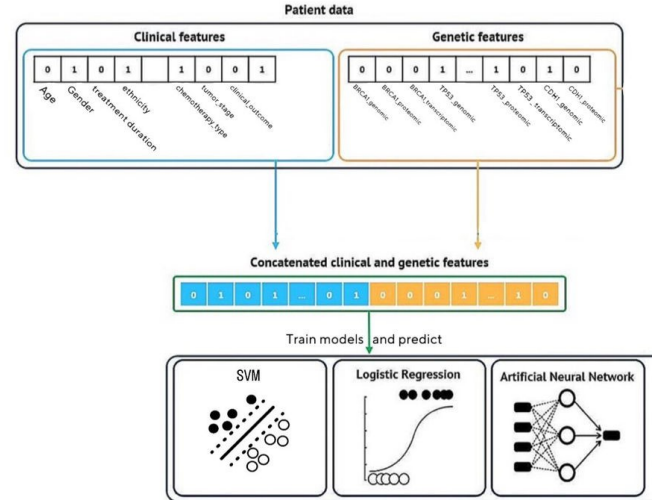


Fig. 1: Overview of the methodology workflow. This figure illustrates the data sources, preprocessing, model training, and evaluation pipeline used in the study.

more interpretable, were used to cross-validate the findings of the MLP, showing consistent patterns in the data

## III. RESULTS

The analysis began with an exploration of the synthetic dataset, which revealed an imbalanced distribution of resistant and non-resistant samples. The dataset included a total of 5,000 samples, with key genomic, transcriptomic, and proteomic features specifically associated with breast cancer chemotherapy resistance. Despite the imbalance, preprocessing steps were taken to improve the quality of the dataset by addressing missing values and standardizing feature scales, ensuring optimal performance of the predictive models

### A. Correlation Analysis

We have performed a broad correlation analysis to investigate the relationships between multi-omics features and resistance to chemotherapy. There were mild correlations of resistance with some genomic features. Resistance positively correlated 0.47 with TP53 genomic expression and negatively with BRCA1 genomic expression, correlating -0.47. This therefore showed that all those variables, though important in chemotherapy resistance, are individually associated in a rather moderate way. In general, the correlation matrix was poor. between features and the resistance label is likely due to the artificial nature of the dataset. Although synthetic data are very useful for controlled experimentation and reproducibility, they cannot capture all the complex relationships and variability that might exist in a real-world multi-omics dataset. This is a limitation that stresses the need for validation on clinical datasets to extend the applicability of findings. Considering these modest correlations, the analysis underlines how such a complex, multifactorial phenomenon as chemotherapy resistance cannot be captured without multi-omics integration.

The machine The learning models developed in this study were designed to account for these subtle patterns by leveraging the combined predictive power of genomic, transcriptomic, and proteomic features.
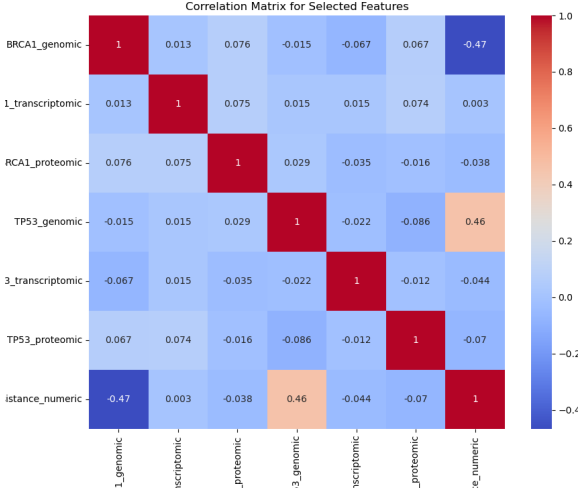


Fig. 2: Correlation matrix of the multi-omics dataset. The matrix shows mild correlations between features and chemotherapy resistance, with a positive correlation for TP53 and a negative correlation for BRCA1, reflecting the synthetic nature of the data.

## B. Performance Metrics

Table I summarizes the performance metrics of the machine learning models evaluated in this study. The MLP demonstrated the highest overall accuracy, precision, recall, and AUC-ROC, outperforming Logistic Regression and SVM across all metrics.

TABLE I: Performance Metrics of Machine Learning Models

| Model | Accuracy | Precision | Recall | AUC-ROC |
|---|---|---|---|---|
| Logistic Regression | 90% | 92% | 76% | 0.96 |
| SVM | 90% | 93% | 75% | 0.96 |
| MLP | 98.51% | 99% | 98% | 1.00 |

## C. Logistic Regression

Logistic Regression achieved a training accuracy of 91.55% and a test accuracy of 90%. As shown in Fig. 3, the confusion matrix revealed that the model performed well in identifying non-resistant cases, with a precision of 92% and recall of 95%. However, its sensitivity to resistant cases was lower, achieving a recall of 76%. This indicates that while Logistic Regression effectively identifies patients likely to respond to chemotherapy, it may underestimate resistance.

## D. Support Vector Machine (SVM)

Support Vector Machine (SVM) delivered comparable results to Logistic Regression, with a test accuracy of 90% and an average cross-validation score of 90%. As depicted
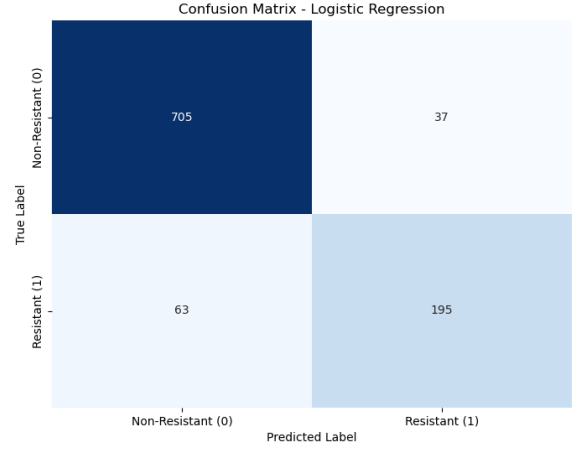


Fig. 3: Confusion matrix for Logistic Regression, showing high performance in identifying non-resistant cases.

in Fig. 4, the SVM demonstrated a balance between precision and recall across both classes, with an AUC-ROC of 0.96. This highlights its robustness in managing the high-dimensional data structure and maintaining predictive stability.
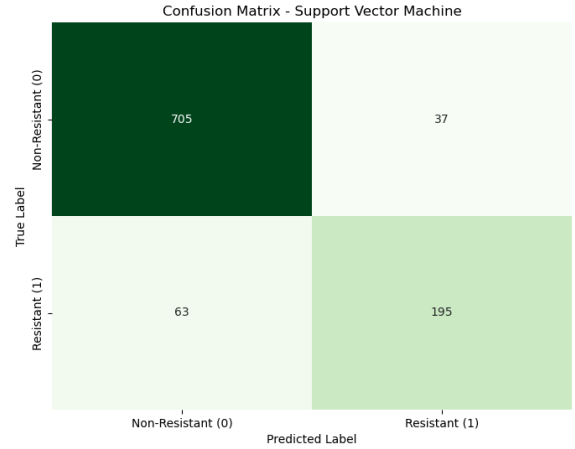


Fig. 4: Confusion matrix for Support Vector Machine (SVM). The SVM exhibited balanced precision and recall across both classes.

## E. Multi-Layer Perceptron (MLP)

Multi-Layer Perceptron (MLP) emerged as the most effective model, achieving a training accuracy of 99.87% and a test accuracy of 98.51%. As shown in Fig. 5, the MLP demonstrated exceptional precision (99%) and recall (98%) across both classes, resulting in a near-perfect AUC-ROC of 1.0. The architecture of the MLP, with two hidden layers and early stopping, enabled it to capture complex, non-linear patterns within the data while avoiding overfitting.
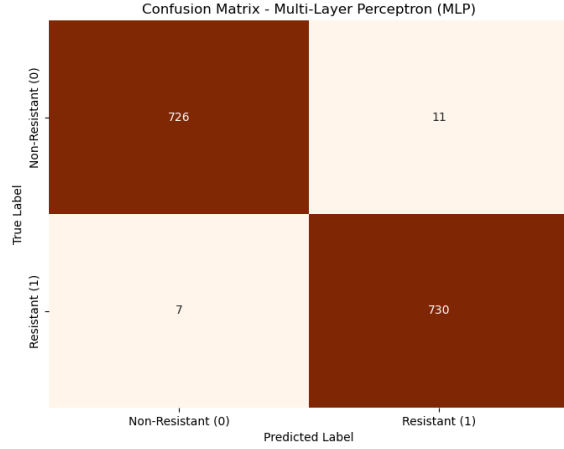
Fig. 5: Confusion matrix for Multi-Layer Perceptron (MLP), showing outstanding classification performance with minimal errors.

## F. Comparison and ROC Curves

The comparison of model performances underscores the strength of deep learning in handling high-dimensional, multi-omics data. However, the consistent results from Logistic Regression and SVM emphasize the utility of simpler models in scenarios where interpretability is prioritized. The ROC curves for all three models, presented in Fig. 6, highlight the superior performance of the MLP, with an AUC-ROC of 1.0 compared to 0.96 for Logistic Regression and SVM. These findings suggest that while advanced architectures like the MLP excel in predictive power, the integration of interpretable models can offer complementary insights, particularly in clinical applications where transparency is essential.
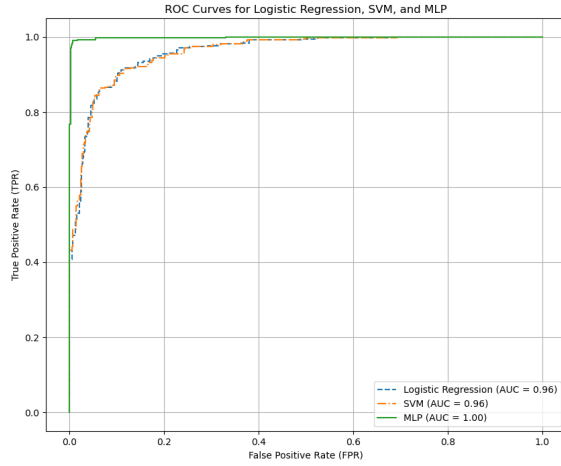


Fig. 6: ROC curves for Logistic Regression, SVM, and MLP. The MLP achieved a perfect AUC of 1.0, outperforming the other models.

## G. Key Observations

The evaluation metrics, confusion matrices, and correlation analysis consistently highlighted the strengths and limitations of each approach. The MLP excelled in distinguishing resistant and non-resistant cases, while Logistic Regression and SVM provided robust baseline performances. The results demonstrate the feasibility of leveraging synthetic multi-omics data to model chemotherapy resistance, paving the way for future applications in personalized oncology.

## IV. DISCUSSION

These results really pinpoint the transformative power of machine learning, especially deep learning models represented by MLP, in predicting chemotherapy resistance using multi-omics data. MLP outperformed Logistic Regression and SVM with near-perfect classification accuracy and very few misclassifications, hence underlining the suitability of deep learning approaches to deal with the inherent complexity and high dimensionality of multi-omics datasets (6, 11, 15). Such performance is essential for real-world clinical applications, where treatment outcomes depend heavily on accurate resistance predictions.

Among the most outstanding results was the high precision of the MLP in nonlinear, high-dimensional data processing, reflected by its outstanding sensitivity and specificity. This capability is particularly critical in clinical applications since treatment plans must be tailored to whether a case is resistant or not. The confusion matrices and ROC curves demonstrated the efficacy of MLP in striking a balance between minimizing false negatives and false positives.While high sensitivity ensures that the resistant cases are picked up as early as possible, high specificity helps avoid the treatments of those patients who are most likely to respond to chemotherapy (6, 13). Such precision goes to directly impact the clinical outcomes because on time detection of resistance has allowed personalization of treatment regimes and evasion of adverse effects due to ineffective therapies.

Synthetic multi-omics data played a pivotal role in this study, offering controlled experimentation and reproducibility while eliminating the ethical and logistical challenges associated with real-world clinical data. Synthetic data enabled consistent conditions for evaluating and comparing models, highlighting the robust performance of the MLP (13, 6). However, synthetic datasets have inherent limitations—they cannot capture the full complexity, variability, and potential confounders present in real-world clinical data (6, 9, 15). The current study, therefore, lays very good ground for the feasibility of machine learning in predicting chemotherapy resistance, but these studies should be complemented in the future by real-world multi-omics data to validate such findings and assure their robustness and clinical applicability.

Deep learning models, including MLP, have been unparalleled in their predictive power, especially for multi-omics applications. However, one of the important clinical barriers to adoption is related to the "black-box" nature of these models. Clinicians often rely on interpretable models that offer

insight into the reasoning for the predictions, such as Logistic Regression or SVM, even if these models are less powerful 16, 17. The challenge is to bridge the gap between deep learning's exceptional predictive accuracy and the interpretability required for real-world decision-making. It can therefore be expected that recent Explainable AI techniques may also make deep learning more transparent and hence more acceptable for the clinicians' interest in 16, 17.

Another major limitation of this study is the use of an imbalanced synthetic dataset. In clinical practice, datasets are often highly imbalanced, with resistant cases being under-represented. While such datasets pose challenges for training and evaluating models, they mirror the real-world challenges more closely than balanced datasets. Advanced techniques, such as synthetic minority over-sampling or cost-sensitive learning, will be required to address this issue and enhance the generalizability of machine learning models in real-world settings (7, 18)

This study focuses not only on the accuracy of the predictive models themselves but also presents synthetic data as a useful resource for preliminary model development. The high recall of MLP is especially remarkable, as it minimizes false negatives; thus, resistant cases will be identified in time to adapt treatment strategies correspondingly. Identifying resistance earlier can significantly help improve patient outcomes by avoiding ineffective treatments and facilitating timely interventions (6, 13). Besides, incorporation of other data modalities, such as imaging or clinical features, might further improve the predictive accuracy and reliability of these models (14, 15).

Future studies should focus on several key aspects. First, real-world multi-omics data should be integrated to validate the findings of this study and ensure robustness. Second, improving the interpretability of deep learning models through explainable AI methods is essential to bridge the gap between predictive performance and clinical trust (16, 17). Finally, addressing the challenges of class imbalance through advanced data augmentation techniques and adaptive learning approaches will be vital for expanding the clinical applicability of these models (7, 18). Overcoming these limitations will bring machine learning closer to reshaping oncology care and offering ever more personalized, effective treatments to patients.

## V. CONCLUSION

This is just another indication that the value of machine learning, mainly deep learning in this regard, may lie in providing predictions with regard to resistance toward chemotherapy in synthetic multi-omics data associated with breast cancer. The model proposed-MLP-scored fantastic and gave an accuracy of 98.51% with 1.0 AUC-ROC during testing. The result outperformed those belonging to both Logistic Regression and SVM. These results show how far deep learning has gone in mastering the complexities from multi-omics data and personal oncology work processes. In addition, synthetic data provided a strong basis for evaluation of models and

then allowed controlled experimentation and reproducibility when it was not possible on real-world datasets. On the other hand, validation in clinical data would be requisite for general applicability and also clinical relevance. Moreover, this lack of interpretability within deep learning algorithms is an outstanding barrier to the clinical integration of these algorithms. Future studies shall focus on real-world data validation of these models, inclusion of multimodal data types, and formulation of explainable AI techniques that make such models more interpretable. Accomplishment of these issues puts machine learning in the best place for significant advances toward personalized oncology and the benefit of better outcomes for the patient

## REFERENCES

[1] World Health Organization. (2024). Breast cancer fact sheet.
[2] World Cancer Research Fund. (2022). Breast cancer statistics.
[3] Pusztai, L., Márk, L., Márk, L. (2016). Chemotherapy in breast cancer: The need for improved biomarkers. *Nature Reviews Clinical Oncology*, *13*(3), 118–130.
[4] Harbeck, N., Gnant, M. (2017). Breast cancer. *The Lancet*, *389*(10085), 1134–1150.
[5] Turner, N., Reis-Filho, J. S. (2018). Genomics and precision medicine for breast cancer treatment. *Cancer Cell*, *34*(5), 709–720.
[6] Vasan, N., Baselga, J., Hyman, D. M. (2019). A view on synthetic data and machine learning in cancer research. *Nature Reviews Drug Discovery*, *18*(6), 383–398.
[7] Wang, L., Meijering, E. (2017). Multi-omics integration for breast cancer classification. *Nature Biotechnology*, *35*(12), 1115–1116.
[8] Sun, W., Wang, H. (2018). Multi-omics analysis of breast cancer subtypes and resistance. *Cell*, *173*(2), 348–360.
[9] Ritchie, M. D., Holzinger, E. R., Li, R., et al. (2015). Integrative analysis of genetic and omics data using machine learning. *Nature Reviews Genetics*, *16*(2), 85–97.
[10] van Dijk, D., Sharma, R. (2021). The role of machine learning in chemotherapy resistance models. *Bioinformatics*, *37*(15), 2076–2085.
[11] Altaf, F., Anwar, S. M., Gul, N. (2021). Deep learning in precision oncology: The challenges and opportunities. *Computers in Biology and Medicine*, *130*, 104183.
[12] Collins, F. S., Varmus, H. (2015). A vision for the future of genomics research in cancer. *New England Journal of Medicine*, *372*(1), 99–107.
[13] Liu, Z., Wu, J., Wang, C. (2020). Synthetic datasets in oncology research: Opportunities and challenges. *Frontiers in Oncology*, *10*, 2051.
[14] Cheng, W. Y., Fuller, R. (2018). Integration of deep learning and clinical features for chemotherapy resistance. *Journal of the American Medical Informatics Association*, *25*(6), 678–685.
[15] Topol, E. J. (2019). High-performance medicine: The convergence of artificial intelligence and healthcare. *Nature Medicine*, *25*(1), 44–56.
[16] Reddy, S., Zhang, L. (2020). The importance of model interpretability in clinical AI applications. *BMJ Innovations*, *6*(3), 92–96.
[17] Venkatesh, S., Suresh, J. (2022). Explainable AI in deep learning for cancer treatment planning. *Artificial Intelligence in Medicine*, *122*, 102191.
[18] Esteva, A., Robicquet, A., Chou, K. (2017). Artificial intelligence for oncology: A primer. *Nature Medicine*, *23*(1), 70–75.