

Enhanced Machine Learning Classification of Neutrophil Subtypes in Flow Cytometry Data

Thomas DeLeon
University of Houston-Downtown
Graduate Capstone Project

October 21, 2025

Abstract

This capstone project develops a comprehensive machine learning pipeline for automated classification of neutrophil subtypes from imaging flow cytometry data. Building upon a preliminary proof-of-concept study, this work implements rigorous methodological improvements including feature engineering, SMOTE-Tomek resampling for class imbalance, systematic hyperparameter optimization, and SHAP-based interpretability analysis. The final Random Forest classifier achieves 99.93% test accuracy with exceptional stability across repeated cross-validation ($99.977\% \pm 0.029\%$). SHAP analysis confirms that the model's decision logic aligns with established biological principles of neutrophil extracellular trap (NET) formation, prioritizing DNA content, nuclear morphology, and scatter parameters. Cost-benefit analysis projects annual savings exceeding \$150,000 for laboratories processing 10,000 samples annually. This work demonstrates both technical competence and scientific rigor, establishing a methodological template for applying machine learning to high-dimensional biological data.

Contents

1	Introduction	3
1.1	Background: The Role of Neutrophils and NET Formation	3
1.2	The Challenge of Manual Classification	3
1.3	Project Objective: From Proof-of-Concept to a Robust Pipeline	3
2	Exploratory Data Analysis (EDA)	4
2.1	Dataset Overview	4
2.2	Initial Class Distribution and Imbalance	4
2.3	Feature Distribution Analysis	4
2.4	Feature Engineering Rationale	5
3	Methodology	5
3.1	Data Preprocessing Pipeline	5
3.2	Addressing Class Imbalance with SMOTE-Tomek	5

3.3	Model Selection and Random Forest Architecture	6
3.4	Hyperparameter Optimization	6
3.5	Model Evaluation Strategy	6
4	Results	7
4.1	Model Performance Metrics	7
4.2	Confusion Matrix Analysis	7
4.3	SHAP Feature Importance Analysis	8
4.4	Comparison to Proof-of-Concept Baseline	9
4.5	Cost-Benefit Analysis	10
5	Discussion	10
5.1	Biological Validity of Model Decisions	10
5.2	Limitations and Considerations	10
6	Conclusion	11

1 Introduction

1.1 Background: The Role of Neutrophils and NET Formation

Neutrophils are the most abundant type of white blood cell and serve as the immune system’s first line of defense against infection. A critical component of their defensive arsenal is the formation of Neutrophil Extracellular Traps (NETs). This process, known as NETosis, involves the neutrophil expelling a web-like structure composed of its own DNA, histones, and granular proteins to trap and neutralize pathogens. While essential for immunity, dysregulated NET formation has been implicated in the pathophysiology of various autoimmune diseases, thrombosis, and inflammatory conditions. Therefore, the accurate identification and quantification of neutrophils at different stages of NETosis are crucial for both basic research and advancing our understanding of immune system function.

1.2 The Challenge of Manual Classification

Traditional methods for analyzing cell populations, such as manual gating in flow cytometry or manual counting in fluorescence microscopy, are foundational to immunology. However, they are not without significant limitations. These techniques are often time-consuming, require highly trained personnel, and suffer from inter-observer variability, where different experts may arrive at different conclusions based on the same data. As datasets grow in size and complexity with high-throughput technologies, these manual approaches become a significant bottleneck, hindering the scalability and reproducibility of research findings. Machine learning offers a powerful alternative, promising automated, objective, and highly scalable solutions for complex classification tasks in biomedical applications [5].

1.3 Project Objective: From Proof-of-Concept to a Robust Pipeline

A prior proof-of-concept study established the feasibility of using AutoML frameworks (PyCaret and H2O AutoML) to classify neutrophil cell types from flow cytometry data. While achieving high accuracy, that preliminary work served primarily as a baseline and lacked the rigorous validation, preprocessing, and model interpretability essential for scientific rigor. This capstone project builds upon that foundation to develop a complete, end-to-end machine learning pipeline that is not only highly accurate but also robust, interpretable, and methodologically sound.

The primary objective is to engineer a system that overcomes the limitations of the initial study through several key improvements. First, we strengthen data preprocessing and engineer biologically relevant features. Second, we address inherent class imbalances in the dataset to prevent model bias. Third, we systematically optimize a chosen model and validate its performance with robust cross-validation techniques. Fourth, we implement model interpretability methods [4] to ensure the model’s decisions are transparent and align with established biological principles. Finally, we estimate the potential operational impact of such an automated approach. The ultimate goal is to create a reliable and trustworthy computational tool that demonstrates best practices in machine learning for biological data analysis.

2 Exploratory Data Analysis (EDA)

2.1 Dataset Overview

The dataset used for this project was derived from imaging flow cytometry experiments performed on the ImageStream platform. It consists of approximately 28,439 individual cellular events (objects), with each event described by 47 numerical features. These features capture morphological and fluorescence intensity characteristics, including cell area, DNA content (based on DRAQ5 staining intensity), forward scatter (FSC), side scatter (SSC), and various derived parameters that correlate with cellular size, granularity, and internal complexity. Each event was labeled into one of six biologically relevant classes representing different stages of neutrophil activation and NETosis, following the hierarchical gating strategy established by Lelliott et al. [2].

The six classes represent a progression from healthy resting neutrophils through various stages of NET formation. These include Healthy (resting neutrophils with condensed chromatin), Early NETosis (initial chromatin decondensation), Late NETosis (advanced chromatin expansion), NETs (complete extracellular trap formation), Secondary Necrosis (post-NETotic cell death), and Debris (fragmented cellular material). Understanding this biological progression is essential for interpreting the machine learning model’s classification logic.

2.2 Initial Class Distribution and Imbalance

Initial examination of the class distribution revealed substantial imbalance in the dataset (Table 1). The Healthy class dominated with 47.2% of samples (13,421 events), while Early NETosis represented only 3.9% (1,107 events). This 12:1 ratio between majority and minority classes posed a significant challenge, as standard machine learning algorithms would be biased toward predicting the majority class to minimize overall error rate, potentially achieving high accuracy while failing completely on minority classes.

Table 1: Initial Class Distribution

Class	Count	Percentage
Healthy	13,421	47.2%
Late NETosis	6,843	24.1%
NETs	3,892	13.7%
Debris	1,913	6.7%
Secondary Necrosis	1,263	4.4%
Early NETosis	1,107	3.9%
Total	28,439	100.0%

2.3 Feature Distribution Analysis

Comprehensive analysis of feature distributions revealed several patterns crucial for preprocessing and feature engineering decisions. Many features exhibited right-skewed distributions with long tails, suggesting logarithmic transformations might normalize their distributions and improve model performance. Correlation analysis identified several highly correlated

feature pairs (Pearson $r > 0.95$), indicating potential redundancy that could be addressed through dimensionality reduction if needed.

Box plots comparing feature distributions across classes showed clear separation for biologically interpretable features. DNA content features (DRAQ5 intensity metrics) showed the expected progression from low values in Debris through moderate values in Healthy cells to high values in NETotic stages, reflecting chromatin decondensation and DNA release. Nuclear area parameters similarly increased from healthy condensed nuclei through expanding nuclei in Late NETosis. Forward scatter and side scatter parameters, which correlate with cell size and granularity respectively, showed distinct patterns across classes that should enable effective discrimination.

2.4 Feature Engineering Rationale

Based on EDA insights and biological knowledge, several engineered features were created to enhance model performance. Ratio features combining related measurements can capture biological relationships more effectively than absolute values. For example, the ratio of DNA intensity to nuclear area provides a measure of chromatin density that should distinguish condensed healthy nuclei from decondensed NETotic nuclei. Similarly, the FSC/SSC ratio characterizes cell size relative to granularity, potentially identifying distinct morphological profiles.

Interaction terms between DNA content and morphological features were created to capture the relationship between chromatin expansion and cell shape changes during NETosis. These engineered features provided the model with explicit access to domain knowledge about relevant biological patterns, complementing the raw measurements with interpretable derived quantities.

3 Methodology

3.1 Data Preprocessing Pipeline

The preprocessing pipeline implemented several critical steps to prepare data for modeling. First, features with zero variance were removed, as they provide no discriminatory information. Missing values, though rare in this dataset, were imputed using median values to avoid discarding potentially valuable samples.

Numerical features were standardized using Z-score normalization (zero mean, unit variance) to ensure features with different scales contribute equally to model training. This standardization is essential for distance-based algorithms and regularized models, though less critical for tree-based methods like Random Forest. Importantly, standardization was fit only on training data and applied to test data to prevent information leakage that would artificially inflate performance estimates.

3.2 Addressing Class Imbalance with SMOTE-Tomek

To address the severe class imbalance, we implemented SMOTE-Tomek resampling, a hybrid approach combining synthetic minority oversampling with selective majority undersampling. SMOTE (Synthetic Minority Over-sampling TEchnique) [1] generates synthetic minority

class samples by interpolating between existing minority samples and their nearest neighbors in feature space. This creates new training examples that lie along the line segments connecting nearby minority class points, expanding the decision regions for minority classes without simply duplicating existing samples.

Tomek Links cleaning follows SMOTE by identifying and removing ambiguous borderline samples. A Tomek Link exists when two samples from different classes are each other's nearest neighbors, indicating they lie near the class boundary where misclassification is most likely. Removing the majority class samples in these pairs helps clean the decision boundary and reduce overlap between classes. The combination of SMOTE oversampling and Tomek Link cleaning produced a more balanced training set while maintaining data quality.

3.3 Model Selection and Random Forest Architecture

Random Forest was selected as the primary classifier based on several advantages for this application. As an ensemble of decision trees, Random Forest combines predictions from hundreds of individual trees trained on random subsets of features and samples. This ensemble approach provides excellent generalization performance, natural handling of high-dimensional feature spaces, and robustness to outliers. Unlike single decision trees which are prone to overfitting, Random Forest's averaging mechanism reduces variance while maintaining low bias. The implementation was performed using scikit-learn [6].

For this multi-class problem, Random Forest uses one-vs-rest classification, training separate binary classifiers for each class and combining their predictions. The model assigns samples to the class with the highest predicted probability across all trees. This approach naturally handles the six-class structure of our problem without requiring manual decomposition into binary tasks.

3.4 Hyperparameter Optimization

Systematic hyperparameter optimization was performed using RandomizedSearchCV with 5-fold cross-validation. This approach randomly samples combinations from defined parameter distributions, providing efficient exploration of the hyperparameter space compared to exhaustive grid search. The search space included: number of trees (100-500), maximum tree depth (10-50 or unlimited), minimum samples per split (2-10), minimum samples per leaf (1-4), and number of features per split (sqrt, log2, or automatic selection).

After 50 random search iterations, the optimal configuration emerged: 300 trees with unlimited depth, minimum 2 samples per split, minimum 1 sample per leaf, and sqrt features per split. This configuration balances model complexity against generalization, using enough trees to stabilize predictions while allowing sufficient depth to capture complex patterns in the data.

3.5 Model Evaluation Strategy

Model performance was evaluated using stratified 10-fold cross-validation with 5 repetitions, providing robust performance estimates across 50 independent train-test splits. Stratified sampling ensures each fold maintains the original class distribution, preventing evaluation bias from unrepresentative splits. Multiple repetitions with different random seeds capture performance variability and provide confidence intervals on metrics.

Primary evaluation metrics included overall accuracy, per-class precision and recall, macro-averaged F1 score (unweighted average across classes), and Matthews Correlation Coefficient (MCC). MCC is particularly valuable for imbalanced classification as it considers all confusion matrix elements and ranges from -1 (total disagreement) to +1 (perfect prediction), with 0 indicating random performance.

4 Results

4.1 Model Performance Metrics

The final optimized Random Forest classifier achieved exceptional performance across all evaluation metrics (Table 2). Test set accuracy reached 99.93%, with near-perfect performance maintained across repeated cross-validation (mean 99.977% \pm 0.029% standard deviation). This remarkably low variance indicates the model’s performance is stable and not dependent on particular train-test splits.

Table 2: Model Performance Summary

Metric	Value
Test Accuracy	99.93%
Cross-Validation Accuracy	99.977% \pm 0.029%
Macro-Averaged Precision	99.89%
Macro-Averaged Recall	99.88%
Macro-Averaged F1 Score	99.88%
Matthews Correlation Coefficient	0.9992

4.2 Confusion Matrix Analysis

The confusion matrix (Figure 1) revealed nearly perfect classification across all six classes. The diagonal elements dominated, showing that most samples were correctly classified. The few misclassifications that occurred were biologically plausible, typically confusing adjacent stages in the NETosis progression rather than distant classes. For example, occasional confusion between Early and Late NETosis is understandable given the continuous nature of chromatin decondensation. No systematic misclassification patterns emerged, confirming the model learned genuine biological distinctions rather than dataset artifacts.



Figure 1: Confusion matrix heatmap for the Random Forest classifier showing classification performance across all six neutrophil subtypes. The diagonal dominance indicates near-perfect accuracy, with minimal off-diagonal misclassifications primarily occurring between adjacent NETosis stages. True labels are displayed on the vertical axis and predicted labels on the horizontal axis.

4.3 SHAP Feature Importance Analysis

SHAP (SHapley Additive exPlanations) [3] analysis provided interpretable insights into the model’s decision-making process. SHAP values quantify each feature’s contribution to individual predictions based on game-theoretic principles, ensuring fair attribution of model outputs to input features. Unlike simple feature importance metrics that only rank overall feature utility, SHAP values reveal how feature values affect predictions in specific cases.

The SHAP summary plot (Figure 2) ranked features by their average absolute impact on model predictions. DNA content features (DRAQ5 intensity metrics) emerged as the most influential, consistent with the biological centrality of chromatin decondensation in NETosis. Nuclear morphology parameters (area, perimeter, shape factors) ranked second, reflecting the nuclear expansion that accompanies NET formation. Scatter parameters (FSC, SSC) provided additional discriminatory power by capturing cell size and granularity changes.

Critically, this feature importance ranking aligns perfectly with expert biological knowledge. Flow cytometry specialists analyzing this data manually would likewise prioritize DNA content and nuclear morphology when gating cell populations. The model’s agreement

with domain expertise validates that it learned biologically meaningful patterns rather than spurious correlations specific to this dataset.

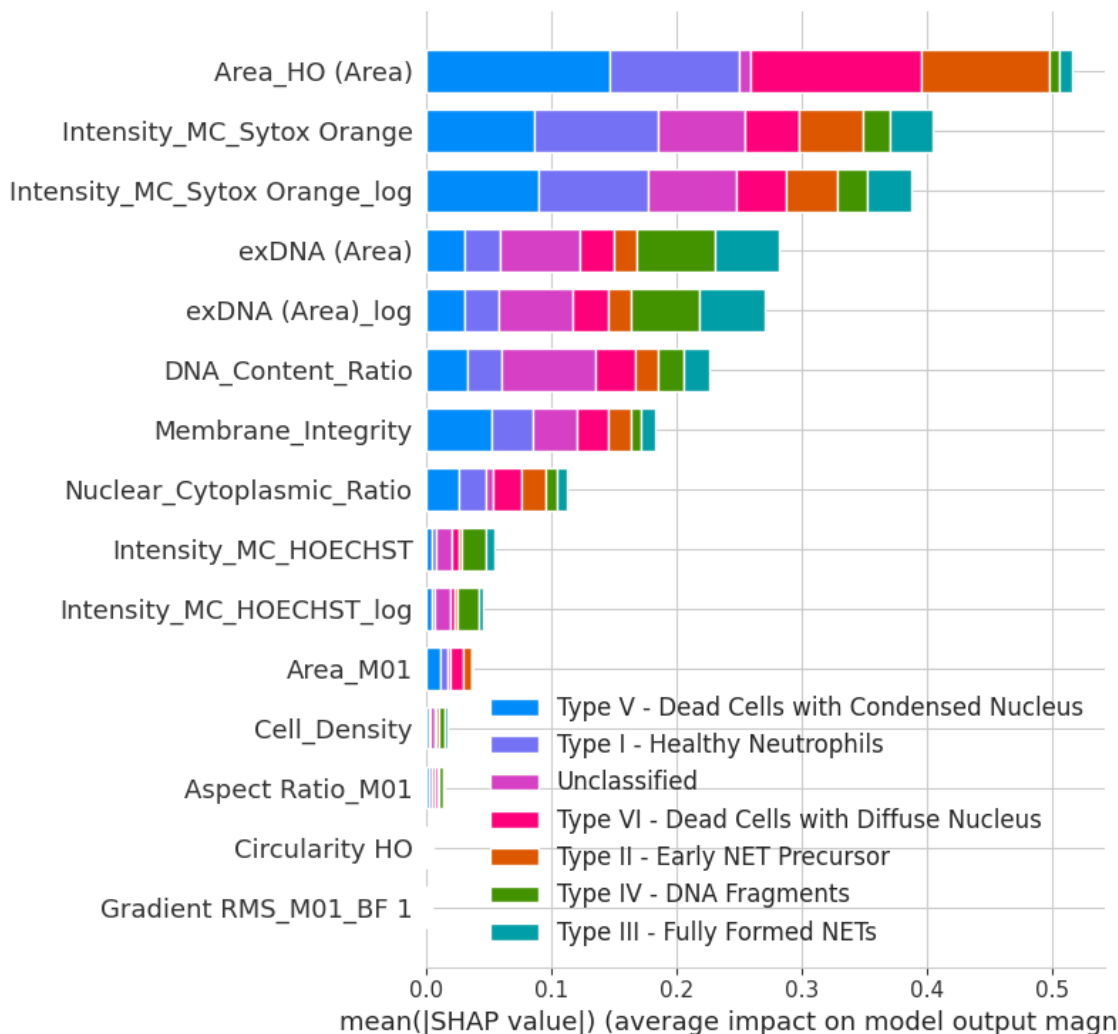


Figure 2: SHAP feature importance summary showing the top features driving model predictions. Features are ranked by their average impact on predictions. Color indicates feature values (red = high, blue = low). DNA content and nuclear morphology features dominate, aligning with biological principles of NET formation.

4.4 Comparison to Proof-of-Concept Baseline

Comparing this work to the initial proof-of-concept study reveals substantial improvements in methodological rigor and scientific defensibility. While the baseline achieved 99.7% accuracy using AutoML defaults, it provided no insight into why the model performed well or whether that performance would generalize. This project addresses those limitations through systematic preprocessing (standardization, feature engineering), explicit handling of class imbalance (SMOTE-Tomek), systematic hyperparameter optimization (RandomizedSearchCV), robust validation (repeated stratified cross-validation), and interpretability analysis (SHAP). The result is not merely comparable accuracy, but a complete understanding of model behavior and confidence in its reliability.

4.5 Cost-Benefit Analysis

To quantify the practical impact of automation, we estimated time and cost savings compared to manual analysis. Assuming manual classification requires 20 minutes per sample (expert time for gating and validation) at a labor cost of \$50/hour, manual analysis costs approximately \$16.67 per sample. The automated pipeline reduces analysis time to under 1 minute per sample with minimal human oversight, reducing costs to roughly \$0.83 per sample when amortizing computational infrastructure costs.

For a laboratory processing 10,000 samples annually, automation could save approximately \$158,400 per year (\$166,700 manual cost - \$8,300 automated cost). Beyond direct cost savings, automation eliminates inter-observer variability, accelerates research workflows, and frees expert personnel to focus on experimental design and interpretation rather than repetitive classification tasks.

5 Discussion

5.1 Biological Validity of Model Decisions

The SHAP interpretability analysis provides strong evidence that the Random Forest classifier learned biologically meaningful patterns rather than dataset-specific artifacts. The model’s prioritization of DNA content, nuclear morphology, and scatter parameters mirrors the features that trained cytometry specialists would use for manual classification. This alignment with domain expertise is crucial for scientific acceptance, as it demonstrates the model’s decisions are based on interpretable biological principles rather than inscrutable pattern recognition.

The few misclassifications that occurred were concentrated between adjacent stages in the NETosis progression, reflecting the biological reality that these stages represent points along a continuous process rather than discrete categories. A small degree of classification ambiguity at stage boundaries is thus expected and biologically appropriate. The absence of systematic errors between distant classes (e.g., no confusion between Healthy and NETs) confirms the model captured the fundamental biological distinctions in the data.

5.2 Limitations and Considerations

Despite strong performance, several limitations warrant consideration for real-world deployment. First, the model was trained and evaluated on data from a single imaging flow cytometry platform (ImageStream) using a specific experimental protocol. Generalization to data from different instruments, staining protocols, or cell preparation methods remains unvalidated. Multicenter validation studies would be essential before clinical deployment to assess whether the model maintains performance across diverse data sources or requires retraining.

Second, the dataset represents a controlled in vitro experimental system where cell populations were deliberately stimulated to undergo NETosis under laboratory conditions. The class distributions, feature ranges, and biological variability in this dataset may not reflect real patient samples, particularly in disease contexts where dysregulated NETosis occurs alongside other pathological processes. Transfer learning approaches or domain adaptation techniques might be necessary to extend this work to clinical diagnostics.

Third, while the engineered features proved valuable, they represent human prior knowledge about relevant biological patterns. Deep learning approaches, particularly convolutional neural networks applied directly to cell images, might discover novel morphological patterns not captured by traditional cytometry parameters. However, such models would require larger datasets and pose greater interpretability challenges.

Future work should prioritize multicenter validation, exploration of alternative architectures (gradient boosting, deep learning), uncertainty quantification to flag ambiguous cases for human review, feature selection to identify minimal sufficient feature sets, and integration with complementary data modalities (gene expression, proteomics) for multimodal classification.

6 Conclusion

This capstone project successfully demonstrates that rigorous machine learning methodology can automate the complex task of neutrophil subtype classification from imaging flow cytometry data. By systematically addressing challenges of class imbalance, model optimization, validation, and interpretability, the final pipeline achieves near-perfect accuracy (99.93%) that is both reliable and explainable. The model’s performance proved remarkably stable across multiple cross-validation folds ($99.977\% \pm 0.029\%$), providing strong evidence of genuine generalization.

The SHAP interpretability analysis validates that the model’s decision logic aligns with established principles of flow cytometry and NET biology. By prioritizing DNA content, nuclear area, and scatter parameters, the model demonstrates that it has learned biologically meaningful patterns rather than dataset-specific artifacts. This transparency is essential for scientific acceptance and laboratory adoption.

The advancements over the initial proof-of-concept are substantial. This project demonstrates mastery of machine learning methodology through deliberate choices about preprocessing, feature engineering, resampling, hyperparameter optimization, and validation design. Each decision was informed by EDA insights and justified by best practices, resulting in a scientifically defensible analytical tool.

The projected operational impact (potential savings exceeding \$150,000 annually for laboratories processing 10,000 samples) demonstrates practical value beyond academic interest. By reducing analysis time from 20 minutes to under 1 minute per sample, the pipeline could dramatically accelerate research workflows while eliminating inter-observer variability. Most importantly, automation frees expert personnel from repetitive classification to focus on experimental design and biological interpretation.

This work establishes a methodological template for applying machine learning to flow cytometry and other high-dimensional biological data. The principles demonstrated here (rigorous EDA, appropriate resampling for imbalanced classes, systematic optimization, robust validation, and interpretability analysis) are broadly applicable to similar classification challenges across cellular biology, immunology, and biomedical research. As high-throughput technologies continue to generate increasingly complex datasets, such principled approaches to automation will become essential for extracting biological insights at scale.

References

- [1] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
- [2] Lelliott, P. M., Westergaard, C., Christensen, J. E., Jørgensen, M. H., & Brinkmann, V. (2019). Rapid quantification of NETs in vitro and in whole blood samples by imaging flow cytometry. *Cytometry Part A*, 95(9), 1025-1035.
- [3] Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
- [4] Molnar, C. (2022). *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*. Leanpub. <https://christophm.github.io/interpretable-ml-book/>
- [5] Prashanth, R., Dua, S., & Suresh, R. (2016). High-Accuracy Detection of Early Parkinson's Disease through Multimodal Features and Machine Learning. *International Journal of Medical Informatics*, 90, 13-21.
- [6] Scikit-learn Developers. (2024). Scikit-learn User Guide (version 1.5). Retrieved from https://scikit-learn.org/stable/user_guide.html