

Prediction of Drug-Induced Autoimmunity Using X Gradient Boost Machine Learning

Srikar Sistla

Dept. of Information Systems
University of Maryland Baltimore County
Baltimore, USA
srikars1@umbc.edu

Kylie Carter

Dept. of Biological Sciences
University of Maryland Baltimore County
Baltimore, USA
kcarter4@umbc.edu

Abstract—Drug-induced autoimmunity (DIA) comprises immune-mediated adverse events such as lupus, hepatitis, and uveitis that can arise after extended drug exposure, complicating prospective risk assessment. We built a gradient-boosted tree (XGBoost) classifier using 196 RDKit-derived molecular descriptors for 477 compounds[1] and addressed class imbalance with SMOTE. On a held-out test set, the model achieved ROC-AUC of 0.888 with 66.7% recall and 57.1% precision for the positive class; five-fold cross-validation indicated strong generalization (ROC-AUC 0.974 \pm 0.067). Gain-based feature importance highlighted topological complexity, aromaticity, and polarity-related descriptors as salient. The framework enables rapid, cost-effective screening of autoimmune risk during early discovery to prioritize compounds for deeper evaluation.

Keywords—Drug-induced autoimmunity, XGBoost, molecular descriptors, SMOTE, preclinical screening, feature importance, gradient boosting, cross-validation

I. INTRODUCTION

Drug-induced autoimmunity (DIA) is an important safety concern in pharmaceutical development. Presentations ranging from drug-induced lupus erythematosus[6] to hepatic[7] or ocular autoimmunity [8] may emerge weeks to months after exposure, leading to missed risks and late-stage setbacks. Heuristic tools such as structural alerts and dose thresholds[5] are informative but often insufficient for prospective triage across novel chemical space.

Structure-based machine learning offers a scalable alternative. Gradient-boosted decision trees (XGBoost) are particularly effective for cheminformatics: they capture nonlinear structure–activity relationships, incorporate regularization to reduce overfitting, and provide descriptor-level interpretability through gain metrics. Here, we develop an XGBoost-based DIA risk model over a curated descriptor set and evaluate it with stratified splitting, SMOTE within training folds, and rigorous cross-validation.

Recent regulatory guidance[4] increasingly emphasizes proactive identification of immunotoxic liabilities earlier in the pipeline to reduce late-stage attrition and post-marketing safety actions. In practice, chemistry teams require lightweight, automatable tools that operate on structure alone, scale to large virtual libraries, and provide rationale for decisions. Descriptor-based gradient boosting fulfills these needs: it integrates seamlessly with cheminformatics workflows, yields probability outputs that can be thresholded for different risk tolerances, and surfaces the most influential features to support discussion with safety, DMPK, and clinical stakeholders.

Another practical challenge is data scarcity and imbalance: autoimmune liabilities are comparatively rare relative to the vast chemical space of benign compounds. This motivates careful validation design, leakage-free resampling (e.g., SMOTE only within

training folds), and attention to calibration so that predicted probabilities meaningfully reflect risk. Beyond point estimates of accuracy, we therefore evaluate confidence via cross-validation and bootstrap intervals, and we discuss threshold selection to tune the precision–recall trade-off for distinct screening scenarios (broad triage versus high-specificity confirmation).

II. LITERATURE REVIEW

A. XGBoost Algorithm:

XGBoost constructs an ensemble by iteratively adding decision trees, where each new tree is trained to minimize the residual errors from the previous trees using gradient descent. The algorithm starts with a simple prediction (often the mean of the target variable) and sequentially adds trees that reduce the loss function. Each tree is constrained by parameters such as maximum depth and minimum child weight to prevent overfitting. The learning rate (0.1 in our case) controls the contribution of each tree to the final prediction, while regularization terms (L1 and L2 penalties) are applied to leaf weights. Subsampling of data points (80% in our configuration) and features further reduces overfitting. The final prediction is the sum of predictions from all trees, weighted by the learning rate. This additive approach allows XGBoost to capture complex, nonlinear relationships while maintaining interpretability through gain-based feature importance, which measures the reduction in loss contributed by each feature across all trees in the ensemble.

B. Current DIA Prediction Methods:

DIA prediction strategies span structural alerts, pharmacokinetic modeling[5], and data-driven learning. Alert and dose heuristics can identify known liabilities but may overcall risk and struggle with unseen scaffolds. Learning-based methods[2] leveraging chemical features have shown improved accuracy and flexibility, especially when paired with transparent model diagnostics.

Emerging hybrid approaches[4] attempt to bridge the gap between mechanistic alerts and data-driven learning by encoding substructure flags, metabolic activation rules, or electrophilicity indices as features within machine learning models. Such methods can capture known causal motifs while allowing the model to learn interactions with physicochemical context (e.g., lipophilicity, polarity, steric bulk). However, reproducibility can suffer when alerts are inconsistently curated across sources, underscoring the value of descriptor-only baselines that are standardized, portable, and easily auditable.

C. XGBoost in Drug Discovery:

XGBoost has become a strong baseline in drug safety modeling due to its bias–variance control, robustness to heterogeneous descriptor scales, and high accuracy across toxicity endpoints. The algorithm’s gain-based importance aids mechanistic interpretation, supporting hypothesis generation while maintaining predictive strength.

Interpretability for gradient boosting has advanced through tools such as SHAP and permutation importance[3], though recent work cautions that explanation methods can be sensitive to correlation structure and feature engineering choices. Consequently, we pair gain-based importance with domain checks (e.g., enrichment of plausible substructures, alignment with known immunogenic motifs) and recommend external validation where chemotype distributions shift. In ADMET practice[4], XGBoost often serves as a strong baseline that can be complemented by conformal prediction for calibrated confidence and by simple rule overlays for conservative screening.

III. METHODOLOGY

Dataset: We used 477 compounds[1] labeled for autoimmune liability (118 positive, 359 negative). Each compound is represented by 196 RDKit descriptors[4] spanning topological indices (e.g., complexity measures), constitutional attributes (e.g., atom and ring counts), and electronic/physicochemical properties (e.g., TPSA, logP).

Data Preprocessing: Features were standardized with StandardScaler. We applied an 80/20 stratified split to preserve class ratios. SMOTE was used only on training folds to mitigate imbalance and avoid leakage into validation/test data.

Model: As explained earlier in Literature Review, we trained an XGBoost classifier with tree learners configured as follows: `n_estimators=300`, `learning_rate=0.1`, `max_depth=6`, `min_child_weight=1`, `subsample=0.8`, with regularization enabled. These settings balance model flexibility and generalization.

Model Evaluation: Metrics included ROC-AUC, precision, recall, and F1-score on the held-out test set. Five-fold cross-validation on the training set assessed stability. Descriptor importance was computed via the gain metric to identify influential features.

IV. RESULTS

A. Model Performance

The XGBoost model achieved a test ROC-AUC of 0.888, indicating strong discrimination. For the positive class, recall was 66.7% and precision 57.1% (F1=0.615), providing a practical balance for early-stage screening where missing liabilities is costly. Cross-validation yielded ROC-AUC 0.974 ± 0.067 , supporting robustness across partitions.

B. Statistical Analysis

Statistical significance testing was performed to validate the model's performance. The 95% confidence intervals for the performance metrics were calculated using bootstrap resampling: ROC-AUC [0.812, 0.938], Precision [0.567, 0.801], Recall [0.425, 0.659], and F1-Score [0.488, 0.722]. These intervals indicate that the model's performance is statistically significant and robust across different data partitions.

C. Feature Importance Analysis:

Gain-based feature importance[3] highlighted descriptors related to topological complexity (e.g., BertzCT), aromaticity (e.g., aromatic ring counts), and polarity (e.g., TPSA). Signals consistent with aromatic amines were among higher-ranked substructure-driven features, aligning with known immunogenic motifs[5]. Importance was distributed across complementary properties, suggesting multi-factor interactions rather than a single dominant cue.

V. COMPARISON WITH LITERATURE

Our XGBoost approach (ROC-AUC: 0.888) demonstrates competitive performance compared to existing DIA prediction methods. [4]

reported ROC-AUC of 0.82 using structural alerts and daily dose information [4], while [2] achieved 0.89 using ensemble methods. Our descriptor-only approach provides comparable performance while retaining interpretability through gain-based feature importance analysis.

The model's precision of 57.1% and recall of 66.7% balance the trade-off between identifying autoimmune-inducing compounds and minimizing false positives, which is crucial for preclinical screening. The cross-validation results (ROC-AUC: 0.974 ± 0.067) indicate robust generalization, suggesting applicability to diverse compound libraries.

VI. CONCLUSION

Our XGBoost model successfully predicts drug-induced autoimmunity with strong performance (ROC-AUC: 0.888) and robust cross-validation (0.974 ± 0.067). The model's recall (66.7%) and precision (57.1%) make it suitable for preclinical screening. Gain-based feature importance identified key molecular descriptors[1], providing mechanistic insights for future validation. This framework offers a practical solution for early-stage DIA prediction[4], potentially reducing late-stage attrition.

REFERENCES

- [1] Huang, X. (2025). Drug Induced Autoimmunity Prediction [Dataset]. UCI Machine Learning Repository.
- [2] Huang, L., Liu, P., & Huang, X. (2025). InterDIA: Interpretable prediction of drug-induced autoimmunity through ensemble machine learning approaches. *Toxicology*, 511, 154064.
- [3] Takefuji Y. (2025). Beyond XGBoost and SHAP: Unveiling true feature importance. *Journal of hazardous materials*, 488, 137382.
- [4] Wu, Y., Zhu, J., Fu, P., Tong, W., Hong, H., & Chen, M. (2021). Machine Learning for Predicting Risk of Drug-Induced Autoimmune Diseases by Structural Alerts and Daily Dose. *International journal of environmental research and public health*, 18(13), 7139.
- [5] Xiao X. and Chang C. (2014) Diagnosis and classification of drug-induced autoimmunity (DIA). *Journal of Autoimmunity*, Academic Press.
- [6] Vaglion A., Grayson P.C., Fenaroli P., Gianfreda D., Boccaletti V., Ghiggeri G.M., Moroni G., (2018) Drug-induced lupus: Traditional and new concepts, *Autoimmunity Reviews*,
- [7] Zeni, N., Cristofani, A., Piano, S. S., Bolognesi, M., & Romano, A. (2025). Pathophysiological Differences and Differential Diagnosis of Autoimmune and Drug-Induced Hepatitis. *Livers*, 5(2), 22.
- [8] Lu, L. M., Wilkinson, V. M., & Niederer, R. L. (2025). Drug-Induced Uveitis: Patterns, Pathogenesis and Clinical Implications. *Clinical Optometry*, 17.