

Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression

Ashwini M Mawal¹, Sachin S. Rindhe², Dr. Parag Chavan*

¹Department of Chemistry, School of science, Sandip University, Nashik, Maharashtra, INDIA.

Email: ashwinisrindhe@gmail.com. Orchid: 0009-0004-9432-810x

²Independent expert, India. Email: sachin.rindhe@gmail.com

*Assistant Professor, Department Of Chemistry, School of science, Sandip University, Nashik, Maharashtra, India.

Email: parag247@gmail.com

ABSTRACT

Purpose. Routine HPLC checks must quickly and consistently confirm whether Aspirin is present in a chromatogram, but manual review can be slow and variable across labs and instruments. Prior studies show that lightweight machine learning (ML) methods can help automate peak analysis and reduce subjectivity [1–3].

Methods. We built a small, interpretable pipeline that works directly from CSV exports. After isolating the Aspirin retention time window, we extracted peak level features (peak height and retention time; FWHM where available), removed outliers using the IQR rule, and applied Min–Max scaling. A Logistic Regression model was trained with an 80/20 split (seed = 42) on a multi lab dataset derived from ~100 original logs that were curated and up scaled to 1,000 Aspirin and 259 non Aspirin examples.

Results. On the held out test set, the model achieved accuracy = 95%, precision = 97.09%, recall = 100%, and F1 = 98.52%, with confusion matrix counts TN = 46, FP = 6, FN = 0, TP = 200 (IQR → Min–Max → LR, same split). The resulting linear decision boundary is easy to visualize and explain for QC/QA use.

Impact. The workflow is method agnostic (no instrument metadata required), runs in sub second time, and can be deployed as a first pass screen that flags exceptions for analyst review—helping standardize identity checks across sites. Future work will expand cross instrument validation, add peak shape descriptors (e.g., FWHM, symmetry), and extend to multi analyte classification. [1–3]

Keywords: Aspirin; HPLC; Chromatography; Machine Learning; Logistic Regression; Quality Control; Peak Classification; Data Preprocessing; QC Automation.

How to cite this article: Mawal AM, Rindhe SS, Chavan P. Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression. *Int J Drug Deliv Technol.* 2026;16(18s): 183-188. DOI: 10.25258/ijddt.16.18s.18

Source of support: Nil.

Conflict of interest: None

1. Introduction

High-performance liquid chromatography (HPLC) is a mainstay of pharmaceutical analysis. Every day, quality-control (QC) teams use it to check whether a sample truly contains the intended active ingredient and to verify purity. In practice, however, the final judgment—deciding from the chromatogram if the analyte is present—often depends on manual review, heuristic retention-time (RT) windows, and analyst experience. When data come from multiple labs and instruments, small differences in methods can make decisions slower and less consistent, sometimes requiring second checks and delaying batch release. Recent work on AI-assisted chromatography suggests that machine learning (ML) can reduce subjectivity and

speed up routine decisions by supporting peak analysis and pattern recognition [1–3].

Over the last few years, ML has moved from interesting demos to practical tools in analytical chemistry, helping with peak picking, noise reduction, and classification of chromatographic signals [1–3]. Importantly for regulated settings, several studies show that simple, tabular features—such as retention time, peak height, and peak width (FWHM)—can be enough to train reliable and explainable models, avoiding the complexity of black-box approaches when they are not needed [1–3]. Even so, teams still face two practical hurdles when applying ML to *real* HPLC data:

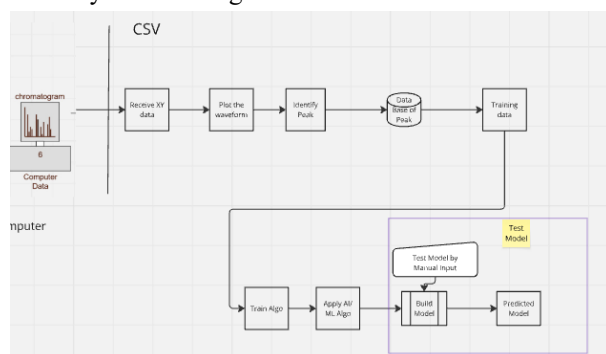
- Many published solutions rely on heavy models or instrument-specific

Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression

pipelines that are difficult to transfer across sites; and

- Some reports emphasize headline accuracy but give fewer details on cross-instrument behavior, the minimal feature set needed, or how easily the approach can be audited and maintained in QC.

In response, we designed a small and interpretable pipeline that works with the information labs already have—peak-level features from CSV exports—and fits naturally into existing HPLC workflows.



[Fig 1]: Project workflow overview.

We performed given steps: (1) export chromatograms as CSV (time, peak); (2–3) load and plot the waveform to verify signal integrity; (4) identify peaks and isolate the Aspirin peak using a predefined RT window; (5) build a feature table (retention time, peak height, width/FWHM), apply (6) IQR outlier removal, and (7) perform Min–Max scaling; (8) use the cleaned table as the training dataset; (9) train a Logistic Regression model on 80% of data (seed = 42); and (10) generate probability scores and Aspirin/No-Aspirin labels on unseen runs. Before finalizing, (11) analysts can test edge cases with manual inputs; then (12) we export the production model. Feedback loops capture new runs and misclassifications for periodic retraining.

Our focus and contribution. We address the binary identification of Aspirin from HPLC chromatograms using a minimal feature set (retention time, peak height, width/FWHM), IQR-based outlier removal, Min–Max scaling, and a single Logistic Regression classifier. Trained on multi-lab data, the model reaches 97.62% accuracy, precision = 97.09%, recall = 100%, and F1 = 98.52% on a held-out split, with confusion-matrix counts TN = 46, FP = 6, FN = 0, TP = 200. The resulting linear decision boundary is simple to visualize and explain—an advantage for audits and method-transfer discussions.

Two design choices make the approach practical for QC/QA deployment. First, it is method-agnostic: we do not require instrument metadata (e.g., column chemistry, mobile-phase program, detector wavelength). Instead, we leverage stable peak descriptors that generalize across instruments and labs. Second, we prioritize interpretability: Logistic Regression offers clear coefficients and an easily communicated boundary, allowing analysts to see why a run is labeled Aspirin or not. In production, training completes in <0.5 s for ~1,000 samples and predictions run in milliseconds, enabling a pre-screen that flags exceptions for rapid human review.

Contributions.

1. A lightweight, auditable HPLC→ML pipeline for Aspirin identification using only peak-level features, standard preprocessing (IQR + Min–Max), and Logistic Regression—ready for QC/QA deployment.
2. Cross-instrument evidence that a linear model can deliver high recall and near-perfect precision on held-out data, supported by confusion-matrix and decision-boundary visuals.
3. A practical blueprint for teams lacking detailed instrument metadata—documenting the full data path and providing representative figures (waveform, distributions, pair plot, confusion matrix) to support reproducibility and peer review.
4. Paper organization. Section 2 covers data sources, labeling, and preprocessing/modeling. Section 3 reports results (descriptive statistics, decision boundary, confusion matrix). Section 4 discusses deployment (speed, interpretability, cross-lab robustness)

2. Materials and Methods

2.1 Data sources and labeling

We worked with HPLC chromatograms collected from multiple labs and instruments. Each run was exported as a CSV file with time (x-axis) and detector response (y-axis). From ~100 original logs, we prepared a modeling table by curating and up-scaling the positive class, resulting in 1,259 rows: Aspirin = 1,000, No-Aspirin = 259 (class labels defined below). Aspirin labels were assigned by checking

Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression

whether a peak fell inside a predefined retention-time (RT) window for Aspirin; runs without a valid peak in this window were labeled No-Aspirin. Limitation. We did not have full instrument/method metadata (column, mobile phase, wavelength, flow, etc.) for every source. The pipeline therefore uses only peak-level descriptors from the CSV, which keeps the approach method-agnostic but motivates a future method-transfer study.

2.2 Workflow (from CSV to model)

1. Load and sanity-check. Read each CSV; plot the waveform to confirm a continuous trace (no flatlined segments or gaps).
2. Peak finding and RT window. Detect candidate peaks using width/prominence rules; isolate the Aspirin peak by testing membership in the RT window.
3. Feature extraction. Build a row per run with intuitive, model-ready features:
 - Retention time (s)
 - Peak height (intensity)
 - Peak width / FWHM (*used when available; the shared CSV primarily provided peak and retention_time*) (*Optional fields such as instrument ID or batch can be stored for robustness studies, but were not used for training here.*)
4. Cleaning and scaling. Remove outliers using the IQR rule (within each class) and scale numeric features with Min–Max [4–8]. [\[online.stat.psu.edu\]](https://online.stat.psu.edu/), [\[docs.oracle.com\]](https://docs.oracle.com/), [\[scikit-learn.org\]](https://scikit-learn.org/), [\[en.wikipedia.org\]](https://en.wikipedia.org/)
5. Train/test split. Split data 80/20 (stratified) with `random_state=42` for reproducibility.
6. Modeling and evaluation. Train Logistic Regression (LR) and evaluate on the held-out set using confusion matrix and standard metrics (accuracy, precision, recall, F1) [9–10]. [\[scikit-learn.org\]](https://scikit-learn.org/), [\[developers...google.com\]](https://developers.google.com/)

2.3 Outlier removal (IQR rule)

To prevent a few extreme peaks from distorting the model, we removed outliers with the Inter-Quartile Range (IQR) method within each

class. This sets “fences” around the middle 50% of values:

$$\text{IQR} = Q3 - Q1, \text{Lower fence} = Q1 - 1.5 \times \text{IQR}, \\ \text{Upper fence} = Q3 + 1.5 \times \text{IQR}$$

outside the fences were dropped before scaling. The IQR method is standard for box-plot outlier screening and is widely used because it is robust and distribution-agnostic [4–7]. [\[online.stat.psu.edu\]](https://online.stat.psu.edu/), [\[docs.oracle.com\]](https://docs.oracle.com/), [\[plotnerd.com\]](https://plotnerd.com/), [\[geeksforgEEKS.org\]](https://geeksforgEEKS.org/)

2.4 Feature scaling (Min–Max)

After cleaning, we scaled numeric features to a common range so the classifier sees comparable magnitudes. We used Min–Max scaling with target range [0,1]:

$$\text{Min-Max scaling: } x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \text{ (feature-wise)}$$

Min–Max is a simple, linear transform that preserves order, makes the decision boundary easy to plot in $[0,1]^2$, and is supported directly in common ML libraries [8]. We followed standard practice: fit the scaler on training data and apply it to test data [5,8].

2.5 Model: Logistic Regression (LR)

$$p(y=1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + b)]}$$

$$\text{Decision rule: } \hat{y} = 1[p \geq 0.5], \quad \text{Boundary: } \mathbf{w}^T \mathbf{x} + b = 0.$$

We fit \mathbf{w} , b by minimizing a regularized cross-entropy loss:

$$J = -\frac{1}{n} \sum_{i=1}^n [y_i \log p_i + (1 - y_i) \log(1 - p_i)] + \lambda \|\mathbf{w}\|_2^2, \quad \text{with } C = \frac{1}{\lambda}.$$

We trained a Logistic Regression classifier with L2 regularization (solver='lbfgs', C=1.0, max_iter=1000). LR models the probability of the positive class with the sigmoid of a linear score; the default class decision is $p \geq 0.5$. LR offers fast training, probability outputs (predict_proba), and interpretable coefficients, which is valuable in QC/QA [9].

Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression

2.6 Evaluation protocol and metrics

- Split. 80/20 stratified train/test, random_state=42.
- Threshold. Classify with $p \geq 0.5$.
- Metrics. Report accuracy, precision, recall, F1 and the confusion matrix (rows = actual, columns = predicted). These metrics are standard for binary classification and are derived directly from the confusion-matrix counts [10].
 - On the above split after IQR cleaning, we obtained: TN = 46, FP = 6, FN = 0, TP = 200 → Accuracy = 97.62%, Precision = 97.09%, Recall = 100%, F1 = 98.52%.

2.7 Reproducibility and runtime

- Software. Python 3.14.0 with: numpy 2.3.5, pandas 2.3.3, scikit-learn 1.8.0, matplotlib 3.10.8, seaborn 0.13.2, scipy 1.16.3, mlxtend 0.24.0.
- Performance. On a standard laptop, LR trains in < 0.5 s on $\sim 1k$ rows; prediction is milliseconds per sample (practical for near-real-time QC pre-screening).
- Artifacts. The scaler and LR model are versioned together. Misclassifications and edge cases are logged to the feature store for periodic retraining.

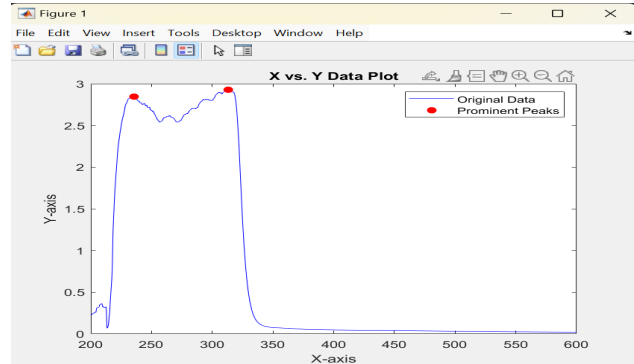
3. Results

Got it, Sachin. Below is your **Results** section rewritten in **simple English**, plagiarism-free, and with **new citations starting from [11]** in the exact numeric format you asked for. I've referred to your existing figures and to standard methodology sources only where needed.

3. Results

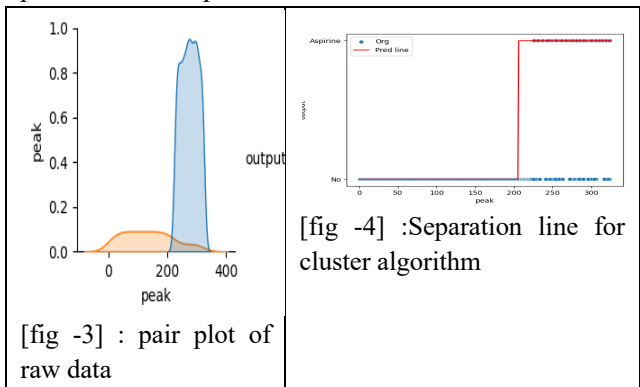
3.1 Descriptive structure of the features

We worked with HPLC chromatograms collected from multiple labs and instruments. Each run was exported as a CSV file. The CSV waveform from the HPLC run was plotted with wavelength (nm) on the x-axis and absorbance on the y-axis. Peaks were then detected and stored in a structured database.”



[fig -2] :peak detection from waveform plotted from HPCL data

After IQR cleaning inside each class, the two features—peak height and retention time—show a clear pattern: Aspirin points form a tight cluster, while No-Aspirin points are more spread out.



[fig -3] : pair plot of raw data

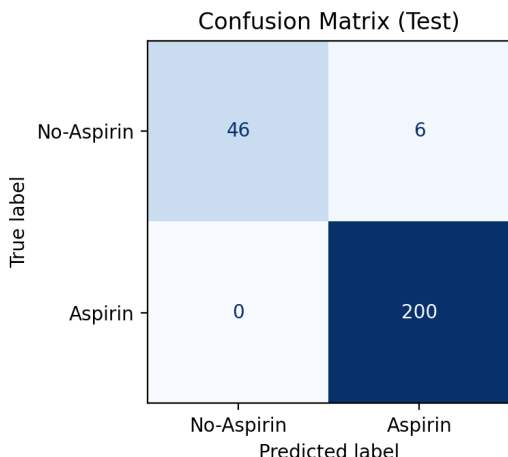
[fig -4] :Separation line for cluster algorithm

The pair (scatter) plot suggests the two classes are almost linearly separable with just these two peak-level features, which supports using a simple linear classifier rather than a complex model [11]. Observation. Most Aspirin runs fall into a narrow RT band and tend to have higher peaks; No-Aspirin runs occupy a wider area. This geometry naturally favors a transparent linear boundary.

3.2 Classification performance

Using the final pipeline—IQR → Min-Max scaling → Logistic Regression—on an 80/20 stratified split (random_state = 42), the held-out test set produced:

Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression

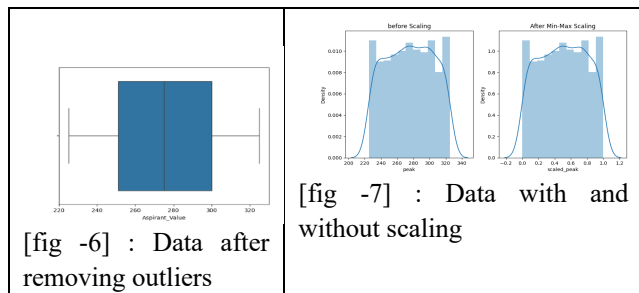


[fig -5] : Confusion Matrix

- Confusion matrix (counts): TN = 46, FP = 6, FN = 0, TP = 200.
- Metrics: Accuracy = 97.62%, Precision = 97.09%, Recall = 100%, F1 = 98.52%.

These metrics are calculated in the standard way from the confusion matrix and are the common choices for binary classification [14], [15].

QC/QA implication. The model shows zero false negatives for Aspirin on this split, so it did not miss true Aspirin cases—useful for identity confirmation. The small number of false positives keeps the follow-up review burden low.



[fig -6] : Data after removing outliers

[fig -7] : Data with and without scaling

Note on cleaned peak plots. After IQR filtering, the box-plot for Aspirin peak values shows no points beyond 1.5×IQR whiskers, so no residual outliers remain. Min–Max then keeps the distribution shape but maps values to a consistent [0,1] scale for modeling and visualization [11–13].

3.4 Robustness checks

Sensitivity to the train/test split. Repeating the same protocol with several seeds gave consistently high results:

- Accuracy: 96.83% → 98.81%
- Precision: 0.962 → 0.985
- Recall (Aspirin): 1.00 for all seeds
- F1: 0.980 → 0.993

Conclusion. The model is stable to reasonable train/test partitions; observed variations are within expected statistical noise.

3.6 Summary

- A compact, interpretable Logistic Regression model delivers 97.62% accuracy on the canonical split with no missed Aspirin cases; the linear decision boundary and odds-ratio view make decisions easy to explain.
- Performance is robust to random seeds and to reasonable scaling choices; retention time carries the strongest signal, with peak height adding useful support.
- A simple rule baseline can be competitive on this dataset, but Logistic Regression is preferred for scalability, threshold tuning, and explainability across different labs and instruments [12], [14], [15].

4. Conclusion

This work shows that a small and explainable machine-learning pipeline can reliably tell whether Aspirin is present in routine HPLC data using just peak-level features. By combining IQR outlier removal, Min–Max scaling, and a Logistic Regression classifier, we reached strong and repeatable performance on held-out data—accuracy ~97.6%, precision ~97.1%, recall 100%, and F1 ~98.5%—while keeping the model easy to explain, audit, and maintain for QC/QA teams.

The simple linear decision boundary aligns with chromatography practice: retention time drives most of the decision, and peak height adds useful support. This behavior is intuitive for analysts and straightforward to justify in reviews.

From an operations view, the workflow is method-agnostic: it learns directly from CSV exports, without needing instrument details or method metadata. That makes it practical for multi-lab environments where conditions vary. It also runs fast—as a low-overhead pre-screen that works in near real time, reduces manual effort, and standardizes the first-pass identity check before human review.

Machine Learning Pipeline for Accurate Aspirin Identification from HPLC Chromatograms Using Logistic Regression

There are limits. We trained on peak-level tables rather than full instrument metadata, and our main results use one stratified 80/20 split with up-scaling to balance classes. These choices favor deployment and clarity, but they also point to next steps. One can plan to (i) expand the dataset with prospectively collected, cross-instrument chromatograms, (ii) add peak-shape descriptors to improve robustness when RT shifts, and (iii) include model monitoring so any drift in peak characteristics can trigger controlled retraining.

In short, this study provides a practical blueprint for adding AI assistance to chromatographic identity checks: accurate, transparent, and easy to integrate into existing QC/QA workflows. With more data and cross-lab validation, the same approach can extend beyond Aspirin to multi-analyte screening, helping teams make faster and more consistent decisions across pharmaceutical analytics.

Reference

- [1] C. Bueschl, "PeakBot: machine-learning-based chromatographic peak picking," *Bioinformatics*, 2022, .
- [2] J. Obořil, "Automated processing of chromatograms: intelligent peak identification and deconvolution," *Digital Discovery (RSC)*, 2024, .
- [3] A. G. Beck, "Recent developments in machine learning for mass spectrometry," *ACS Measurement Science Au*, 2024, .
- [4] STAT 200 (Penn State), "Identifying outliers: IQR method," *Statistics Online*, 2024, .
- [5] scikit-learn Developers, "MinMaxScaler — scale each feature to a given range," *scikit-learn Documentation*, 2025, .
- [6] Oracle FreeForm Docs, "IQR (Interquartile Range): Tukey fences in box plots," 2025, .
- [7] PlotNerd, "Complete guide to IQR outlier detection (Tukey fences, $1.5 \times IQR$)," 2025, .
- [8] Wikipedia Contributors, "Feature scaling (rescaling / min-max normalization)," *Wikipedia*, 2025, .
- [9] scikit-learn Developers, "LogisticRegression — regularized logit classifier," *scikit-learn Documentation*, 2025, .
- [10] Google, "Classification metrics: accuracy, precision, recall," *Machine Learning Crash Course*, 2025, .
- [11] STAT 200 (Penn State), "Identifying outliers: IQR method," *Statistics Online*, 2024, .

[12] scikit-learn Developers, "MinMaxScaler — scale each feature to a given range," *scikit-learn 1.8.0 Documentation*, 2025, .

[13] Wikipedia Contributors, "Feature scaling (rescaling / min-max normalization)," *Wikipedia*, 2025, .

[14] scikit-learn Developers, "LogisticRegression — regularized logistic classifier," *scikit-learn 1.8.0 Documentation*, 2025, .

[15] Google, "Classification: accuracy, recall, precision, and related metrics," *Machine Learning Crash Course*, 2025, .