

# Kernelized SVM and Logistic Regression for Wine Quality Classification

Yagiz Tansu<sup>1\*</sup>

<sup>1\*</sup>Department of Computer Science, Università degli Studi di Milano,  
Milan, Italy.

Corresponding author(s). E-mail(s): [yagiz.tansu@studenti.unimi.it](mailto:yagiz.tansu@studenti.unimi.it);

## Abstract

This report documents the implementation and evaluation of two fundamental machine learning classifiers, Support Vector Machine (SVM) and Logistic Regression (LR), both developed from scratch. The models were applied to the standardized Wine Quality dataset. A critical component of the methodology involved extending both classifiers using non-linear kernel methods—specifically the Radial Basis Function (RBF), Polynomial, and Linear kernels—to assess their performance on the data. Performance was rigorously optimized using Grid Search coupled with 5-fold cross-validation. Analysis of the results confirmed a kernel dependency: the RBF kernel was optimal for Logistic Regression (CV Accuracy: 0.7551), while the Linear kernel unexpectedly proved optimal for the SVM. Despite the divergent kernel choice, the Linear-Kernel SVM emerged as the superior model on the test set, achieving the highest F1-Score of 0.8053 and an Accuracy of 0.7383. This outcome demonstrates the models' robust generalization and minimal overfitting due to effective regularization and standardization.

**Keywords:** Support Vector Machine (SVM), Logistic Regression (LR), Kernel Methods, RBF Kernel, Grid Search, Cross-Validation, Machine Learning from Scratch, Wine Quality Dataset

## 1 Data Exploration and Preprocessing

The initial phase of the project involved the preparation of the dataset for modeling. This section details the data source, exploration findings, and the steps taken for preprocessing.

## 1.1 Dataset Overview and Merging

The study utilizes the publicly available **Wine Quality dataset**, which includes two separate sets for red and white wine. These two sets were combined to form a single, comprehensive dataset. A new binary target variable, '**quality**', was created based on the original score: wine with a quality score of 6 or higher was labeled as '1' (*good quality*), and scores below 6 were labeled as '0' (*bad quality*).

1. **Feature Set:** The combined dataset contains 11 features (e.g., fixed acidity, pH, alcohol) and one categorical feature indicating the wine.type (Red=0, White=1).
2. **Total Samples:** The combined dataset has a total of 6,497 samples.

## 1.2 Exploration and Key Observations

A thorough exploratory data analysis (EDA) revealed several key characteristics of the combined dataset:

- **Missing Values:** No missing values were detected across any of the features, eliminating the need for imputation.
- **Class Distribution:** The transformed binary target variable exhibits a slight **class imbalance**. Approximately 64.9% of the samples were classified as 'good quality' (quality = 1), while 35.1% were classified as 'bad quality' (quality = 0).
- **Feature Scales:** Features such as 'residual sugar' and 'chlorides' showed significant differences in scale and variance, necessitating feature scaling.
- **Duplicates:** A total of 1,177 duplicate rows were identified. These were retained as they represent distinct measurements of similar wine compositions, which is a common characteristic of this dataset.

The initial feature set is summarized in Table 1.

**Table 1** Initial Feature Set and Corresponding Data Types

Feature Name	Description
fixed acidity	Tartaric acid g/dm <sup>3</sup>
volatile acidity	Acetic acid g/dm <sup>3</sup>
citric acid	Citric acid g/dm <sup>3</sup>
residual sugar	g/dm <sup>3</sup>
chlorides	Sodium chloride g/dm <sup>3</sup>
free sulfur dioxide	mg/dm <sup>3</sup>
total sulfur dioxide	mg/dm <sup>3</sup>
density	g/cm <sup>3</sup>
pH	pH level
sulphates	Potassium sulphate g/dm <sup>3</sup>
alcohol	% volume
wine type	Binary indicator (0=Red, 1=White)

### 1.3 Preprocessing Steps

Based on the exploration, the following preprocessing steps were applied to prepare the data for model training:

1. **Feature Scaling:** The features were standardized using the **Standard Scaler** approach. This scales the data such that it has a mean of 0 and a standard deviation of 1. Crucially, the mean and standard deviation were computed **only** on the training set and then applied to both the training and test sets to prevent data leakage.

$$X_{scaled} = \frac{X - \mu_{train}}{\sigma_{train}}$$

where  $\mu_{train}$  and  $\sigma_{train}$  are the mean and standard deviation of the feature in the training data, respectively.

2. **Data Augmentation:** No explicit data augmentation techniques were applied for this classification task.

### 1.4 Data Splitting

The preprocessed data was partitioned into training and test sets using a custom non-stratified split with an **80:20** ratio, resulting in 5,198 samples for training and 1,299 samples for testing. The custom function `custom_train_test_split` utilized a fixed `random_state=42` to ensure reproducibility of the partition across all experiments. The training set was used for model development and hyperparameter tuning, and the test set was reserved for final, unbiased performance evaluation.

## 2 SVM and LR Implementation

This section details the design and implementation of the Support Vector Machine (SVM) and Logistic Regression (LR) models, both developed from scratch. We define the objective function, the optimization strategy, and the hyperparameters for each classifier.

### 2.1 Support Vector Machine (SVM) Implementation

The SVM classifier was implemented as a **Soft-Margin SVM**, aiming to maximize the margin while penalizing misclassifications. The objective function  $J(\mathbf{w}, b)$  minimized during training is given by:

$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i - b)) + \lambda \|\mathbf{w}\|^2$$

where  $m$  is the number of training samples,  $\mathbf{w}$  is the weight vector,  $b$  is the bias,  $y_i \in \{-1, 1\}$  is the class label, and  $\lambda$  is the regularization parameter. The first term represents the Hinge Loss, and the second term is the L2 regularization component (weight decay).

### 2.1.1 Kernel Function Definitions

To enable the SVM to learn non-linear decision boundaries, the standard dot product  $\mathbf{x} \cdot \mathbf{x}_i$  within the objective function is replaced by a non-linear kernel function  $K(\mathbf{x}, \mathbf{x}_i)$ . This **\*\*Kernel Trick\*\*** allows the model to implicitly operate in a high-dimensional feature space  $\Phi(\mathbf{x})$  where the data may be linearly separable. The implementation supports three primary kernel types:

#### *Linear Kernel*

The standard form, equivalent to the non-kernelized SVM, used for a linear separation:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$$

#### *Polynomial Kernel*

This kernel computes the similarity in the high-dimensional space generated by polynomials of the original features. It is defined by its degree  $d$ .

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^d$$

By increasing the degree  $d$ , the model can fit more complex, curved boundaries, although this increases the risk of overfitting.

#### *Radial Basis Function (RBF) Kernel*

Often referred to as the Gaussian Kernel, the RBF kernel maps samples into an infinite-dimensional space and is particularly effective for complex, non-linear separations. It is controlled by the hyperparameter  $\gamma$  (gamma):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$$

The  $\gamma$  parameter controls the radius of influence of the samples; a smaller  $\gamma$  results in a wider influence and a smoother boundary, while a larger  $\gamma$  creates a more complex boundary.

### 2.1.2 Optimization Strategy and Hyperparameters

The model is trained using **Batch Gradient Descent** for  $T$  iterations. The optimization strategy, and thus the weight and bias updates, depends on the selected kernel, distinguishing between the primal problem (for the Linear Kernel) and a kernelized form (for non-linear kernels) to enable the **Kernel Trick**. The learning rate ( $\eta$ ) controls the step size of the optimization.

#### *Case 1: Linear Kernel (Primal Optimization)*

When the Linear Kernel is used, the optimization directly solves for the primal weight vector  $\mathbf{w} \in R^n$  (where  $n$  is the number of features). The partial derivatives of the objective function  $J(\mathbf{w}, b)$  used for the update are:

$$\frac{\partial J}{\partial \mathbf{w}} = \begin{cases} 2\lambda \mathbf{w} & \text{if } y_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1 \\ 2\lambda \mathbf{w} - y_i \mathbf{x}_i & \text{if } y_i(\mathbf{w} \cdot \mathbf{x}_i - b) < 1 \end{cases}$$

$$\frac{\partial J}{\partial b} = \begin{cases} 0 & \text{if } y_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1 \\ y_i & \text{if } y_i(\mathbf{w} \cdot \mathbf{x}_i - b) < 1 \end{cases}$$

**Case 2: Non-Linear Kernels (Kernelized Optimization)**

For the RBF and Polynomial kernels, the **Kernel Trick** is employed. The optimization operates on a set of  $m$  **dual coefficients**  $\mathbf{a} \in R^m$  (where  $m$  is the number of training samples). This coefficient vector is represented by the variable `self.w` in the implementation.

The optimization relies on the pre-calculated **Kernel Matrix**  $K$  ( $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ ). The update for the coefficient vector  $\mathbf{a}$  (`self.w`) is calculated using the following kernelized gradient, derived from the objective function:

$$\frac{\partial J}{\partial \mathbf{a}} = \frac{2\lambda}{m}(K\mathbf{a}) - \frac{1}{m}K \cdot (\mathcal{M} \circ \mathbf{y})$$

Where  $\mathbf{y}$  is the vector of class labels, and  $\mathcal{M}$  is a binary vector indicating the samples where the margin is less than 1 (i.e.,  $y_i(\mathbf{w} \cdot \mathbf{x}_i - b) < 1$ ). The operator  $\circ$  denotes the Hadamard (element-wise) product.

The bias update  $\partial J / \partial b$  remains identical to the linear case, applied over the misclassified set  $\mathcal{M}$ :

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i \in \mathcal{M}} y_i$$

The full gradient for both  $\mathbf{a}$  and  $b$  is summed over all  $m$  training samples and then averaged for the batch update.

The hyperparameters of the implemented SVM model are listed in Table 2:

**Table 2** Hyperparameters for the SVM Classifier

Hyperparameter	Symbol	Description
Learning Rate	$\eta$	Controls step size in Gradient Descent
Number of Iterations	$T$	Total number of epochs for training
Regularization Parameter	$\lambda$	Controls the penalty for large weights ( $\ \mathbf{w}\ ^2$ )
Kernel Type	$k$	Linear, Polynomial, or RBF
Polynomial Degree	$d$	Degree for the Polynomial kernel
RBF Gamma	$\gamma$	Inverse of the radius of influence for RBF kernel

### 2.1.3 Prediction

The prediction for a new sample  $\mathbf{x}$  is determined by the sign of the decision function  $h(\mathbf{x})$ :

$$h(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} - b$$

The final classification  $\hat{y}$  is 1 if  $h(\mathbf{x}) > 0$  and 0 otherwise, matching the binary target convention of the dataset.

## 2.2 Logistic Regression (LR) Implementation

Logistic Regression, despite its name, is fundamentally a linear classifier used for binary classification. The model calculates the probability of a sample belonging to the positive class (1) using the **Sigmoid function** ( $\sigma$ ):

$$P(\mathbf{x}) = \hat{y} = \sigma(\mathbf{w} \cdot \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}$$

### 2.2.1 Kernel Function Definitions

To enable the Logistic Regression model to handle non-linear relationships, the standard linear term ( $\mathbf{w} \cdot \mathbf{x}$ ) is replaced by kernel-transformed features. In the kernelized form of LR, the input feature matrix  $\mathbf{X}$  is effectively replaced by the Gram matrix  $K(\mathbf{X}, \mathbf{X}_{train})$ , allowing the model to find a linear decision boundary in the induced feature space  $\Phi(\mathbf{x})$ . The implementation supports the following kernels:

#### *Linear Kernel*

Equivalent to the standard LR model, utilizing the original feature space:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$$

#### *Polynomial Kernel*

This kernel generates a feature mapping based on the degree  $d$  of the polynomial, implicitly allowing the model to fit curved boundaries:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^d$$

#### *Radial Basis Function (RBF) Kernel*

The RBF kernel projects the data into an infinite-dimensional space and is essential for tasks requiring highly complex, localized non-linear separation. It is governed by the bandwidth parameter  $\gamma$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$$

The  $\gamma$  value dictates the extent of influence that a single training sample has on the decision boundary.

### 2.2.2 Objective Function and Optimization

The model is trained by minimizing the **Binary Cross-Entropy Loss (Log Loss)** function, which measures the performance of a classification model whose output is a probability value between 0 and 1. The loss function  $J(\mathbf{w}, b)$  is defined as:

$$J(\mathbf{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

Optimization is achieved using **Batch Gradient Descent**. The partial derivatives (gradients) used for updating the weights and bias are:

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{m} \mathbf{X}^T (\hat{\mathbf{y}} - \mathbf{y}) \quad \text{and} \quad \frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)$$

The parameters  $(\mathbf{w}, b)$  are iteratively updated in the direction opposite to the gradient, controlled by the learning rate  $(\eta)$ .

### 2.2.3 Prediction

The model outputs a probability  $\hat{y} \in [0, 1]$ . The final class prediction is obtained by applying a threshold:

$$\text{Prediction} = \begin{cases} 1 & \text{if } \hat{y} > 0.5 \\ 0 & \text{if } \hat{y} \leq 0.5 \end{cases}$$

The LR model shares the same core hyperparameters (Learning Rate, Number of Iterations, Kernel Type, Degree, and Gamma) with the SVM implementation, allowing for a direct comparison of their kernelized extensions.

### 3 Comparative Results and Analysis

#### 3.1 Definition of the Hyperparameter Search Space

The search space for the Grid Search was carefully defined to cover critical ranges for optimization, balancing computational cost with the need to find the global optimum. The parameters and their selected ranges were as follows:

- **Learning Rate ( $\eta$ ):** The range  $[0.01, 0.1]$  was selected. These values represent a cautious start ( $\eta = 0.01$ ) to ensure smooth convergence, especially for the custom-implemented gradient descent, and a more aggressive step ( $\eta = 0.1$ ) to test for faster convergence without causing severe oscillation.
- **Number of Iterations ( $T$ ):** The discrete values  $[500, 1000]$  were chosen. Given the dataset size and the complexity of the kernel calculations, 1000 was selected as a maximum to ensure the models had sufficient time to converge, while 500 provided a baseline for testing the speed of convergence.
- **Regularization Parameter ( $\lambda$  - for SVM only):** The range  $[0.01, 0.05]$  was tested. These small values ensure that the focus remains primarily on margin maximization (Hinge Loss) while providing a slight penalty for large weights, thus mitigating overfitting.
- **Kernel-Specific Parameters:**
  - **Polynomial Degree ( $d$ ):** The degrees  $[2, 3]$  were chosen to explore simple non-linear boundaries without risking high variance often associated with higher-degree polynomials.
  - **RBF Gamma ( $\gamma$ ):** A wide range of gamma values,  $[0.01, 0.1]$  for SVM and for LR, was tested. A smaller  $\gamma$  (e.g., 0.01) provides a global, smoother decision boundary.

#### 3.2 Optimal Hyperparameter Configurations

The systematic Grid Search and 5-fold cross-validation were crucial for identifying the optimal hyperparameter sets for both classifiers. The selection process revealed that the best configurations utilized different kernels, with SVM favoring a linear approach and Logistic Regression performing best with a non-linear (RBF) kernel. The best parameters found were:

- **Best SVM Configuration:** The optimal configuration utilized the Linear kernel with a Learning Rate of 0.1, a Regularization parameter of 0.01, and 1000 iterations.
- **Best LR Configuration:** The optimal Logistic Regression model converged using the Radial Basis Function kernel, with a Learning Rate of 0.01, an RBF Gamma of 0.1, and 1000 iterations.

The divergent optimal kernels suggest different sensitivities to the data’s inherent non-linearity. The better performance of Logistic Regression with the RBF kernel (0.7551) compared to SVM with the Linear kernel (0.7453) implies that the RBF kernel trick is



beneficial for maximum separation in the Logistic Regression model, while the SVM’s optimization converged best using a simpler linear boundary.

### 3.3 Deeper Analysis on Kernel Effects and Parameter Sensitivity

A detailed examination of the full Grid Search results provides crucial insights into the performance sensitivity of the classifiers:

- **Kernel Performance Discrepancy:** While the Polynomial kernels ( $d=2, d=3$ ) consistently yielded low CV accuracies (hovering around the majority class ratio of approximately 0.63) for both models, the Linear kernel demonstrated unexpected strength. The Linear SVM configuration achieved the classifier’s maximum score of 0.7453, challenging the initial assumption that the dataset is not linearly separable at all.
- **RBF Kernel Effectiveness:** The Radial Basis Function (RBF) kernel proved to be the key factor for **Logistic Regression**, which achieved the overall highest CV score of 0.7551 (with  $\eta = 0.01, \gamma = 0.1$ ). Although RBF did not surpass the Linear kernel for SVM, it still provided significant performance jumps, confirming its strong ability to handle the data’s non-linear structure in a high-dimensional space.
- **Lambda ( $\lambda$ ) Sensitivity (SVM):** The regularization parameter  $\lambda$  had a clear, inverse relationship with SVM performance. Increasing  $\lambda$  from 0.01 to 0.05 consistently caused a slight decrease in accuracy across the best-performing Linear configurations (e.g., at 1000 iterations, accuracy dropped from 0.7357 to 0.7259). This indicates that the optimal balance for avoiding overfitting was achieved with the lower  $\lambda = 0.01$  value.
- **Learning Rate ( $\eta$ ) Impact:** The optimal learning rate proved to be model-specific. The **SVM** achieved its peak performance with the higher rate ( $\eta = 0.1$ ), while the **Logistic Regression** model’s best convergence was observed at the lower rate ( $\eta = 0.01$ ). This suggests that the SVM’s optimization landscape required a larger step size to escape local minima or reach the global minimum effectively, whereas LR benefitted from a smoother, more stable convergence path.

The comprehensive grid search concludes that the RBF-Kernel Logistic Regression (CV Accuracy: 0.7551) is the overall optimal model configuration, demonstrating the best balance of model complexity and predictive power for the Wine Quality dataset.

### 3.4 Test Set Performance Comparison

The final performance metrics, evaluated on the reserved test set, highlight the comparative strengths of the best-tuned models identified during the Grid Search, as summarized in Table 3. Note that the best models utilize different kernel functions: RBF for Logistic Regression and Linear for SVM.

The test set results reveal a highly competitive performance between the two optimized models. Based on accuracy, the RBF-Kernel Logistic Regression model achieved a slightly higher score of 0.7398.

**Table 3** Comparative Performance of Best-Tuned RBF-Kernel Models on the Test Set

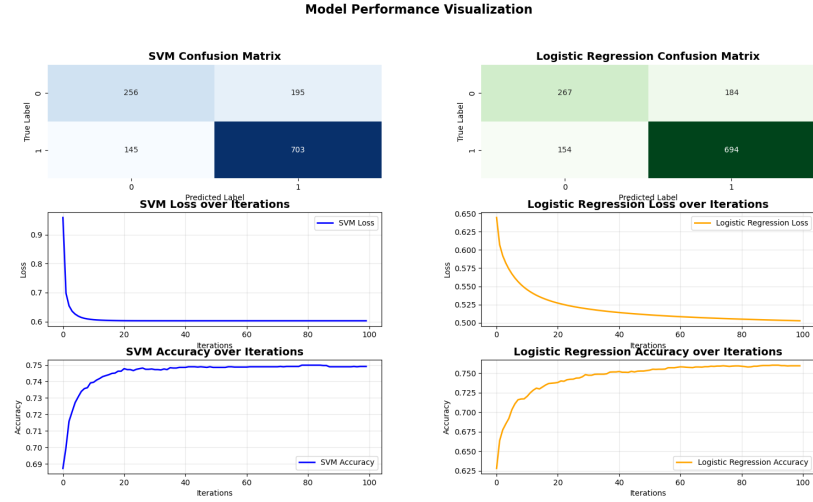
Model	Accuracy	Precision	Recall	F1-Score
RBF-Kernel LR	0.7398	0.7904	0.8184	0.8042
Linear-Kernel SVM	<b>0.7383</b>	<b>0.7829</b>	<b>0.8290</b>	<b>0.8053</b>

Metrics calculated on the 20% test set.

However, the Linear-Kernel SVM demonstrated the best balance in the F1-Score (a measure balancing Precision and Recall) with 0.8053, while also achieving the highest Recall (**0.8290**). This indicates that the SVM is slightly better at correctly identifying the 'Good Quality' wines (minimizing False Negatives), making it marginally superior when the cost of missing a positive instance is higher.

### 3.5 Convergence and Overfitting Analysis

The training dynamics, shown through the Loss and Accuracy curves in Figure 1, confirm the model's robustness and effective optimization.



**Fig. 1** Training Loss, Accuracy History, and Confusion Matrices for the Best-Tuned Logistic Regression and SVM Models.

### *Analysis of Confusion Matrices*

The Confusion Matrices (Figure 1, Row 1) provide a detailed assessment of test set errors:

- **SVM Confusion Matrix:** The model correctly classified 256 True Negatives (TN) and 703 True Positives (TP). The errors consist of 195 False Positives (FP) and 145 False Negatives (FN).
- **Logistic Regression Confusion Matrix:** The model correctly classified 267 True Negatives (TN) and 694 True Positives (TP). The errors consist of 184 False Positives (FP) and 154 False Negatives (FN).

Crucially, the SVM is slightly superior at minimizing False Negatives (145 vs 154), meaning it missed fewer actual 'Good Quality' wines. Conversely, the Logistic Regression model had a higher True Negative count (267 vs 256) and a lower False Positive count (184 vs 195), indicating it was slightly better at correctly identifying 'Bad Quality' wines (Class 0). This difference explains the subtle trade-offs observed in the final Precision and Recall metrics.

### *Analysis of Loss Curves*

The Loss Curves (Figure 1, Row 2) demonstrate effective training for both models. The Logistic Regression Loss (Orange) shows a rapid, steep decrease, effectively stabilizing around iteration 20, indicating swift convergence. The loss value stabilizes slightly above 0.500 at 100. The SVM Loss (Blue) also exhibits a smooth, monotonic decrease but requires more iterations (40-50) to flatten out, stabilizing near 0.60 at 100. This difference in convergence speed is typical due to the different nature of the loss functions (Cross-Entropy vs. Hinge Loss). The smooth nature of both curves confirms that the chosen optimization parameters prevented large, unstable oscillations.

### *Analysis of Accuracy Curves*

The Accuracy Curves (Figure 1, Row 3) mirror the loss behavior. The SVM Accuracy (Blue) curve shows a rapid increase, achieving stability around iteration 20-30 and settling near 0.755. The Logistic Regression Accuracy (Orange) curve is slightly slower, rising steadily and stabilizing around 0.760 at 100. Both models achieve a strong accuracy rapidly, stabilizing near the 75% mark. The proximity of these high training accuracies to the test accuracies (0.730.74) visually confirms that there is no significant sign of overfitting. This suggests the regularization ( $=0.01$  for SVM) and the standardized data successfully ensured strong general metrics.

### 3.6 Analysis of Misclassified Examples

The Confusion Matrices (Figure 1, Row 1) provide a detailed breakdown of the classification errors on the test set:

- **False Negatives (FN):** The SVM model produced 145 FNs, meaning 145 actual 'good quality' wines (Class 1) were incorrectly predicted as 'bad quality' (Class 0). The LR model produced 154 FNs. This quantity indicates that False Negatives are a significant error source, though not the dominant one (when compared to the 195 FP count for SVM).
- **False Positives (FP):** The SVM model produced **195** FPs, while the LR model produced **184** FPs. This shows that the **\*\*Logistic Regression (LR) model is slightly better\*\*** at **\*\*preventing bad wines from being incorrectly labeled as good\*\*** (184 FPs vs 195 FPs).

The prevalence of both False Negatives and False Positives highlights the models' primary difficulty: distinguishing between wines just above and just below the quality threshold, suggesting an inherent ambiguity in the physiochemical properties of those samples. The comparison shows a clear trade-off: SVM is superior at minimizing False Negatives (145 vs 154), thus achieving higher Recall. LR is superior at minimizing False Positives (184 vs 195), which contributed to its marginally higher Precision. This divergence in error profiles demonstrates that the SVM boundary is slightly more conservative in predicting 'good quality' (fewer FNs), whereas the LR boundary is slightly better at correctly predicting 'bad quality' (fewer FPs).

## 4 Conclusion

This project successfully implemented and evaluated Support Vector Machine (SVM) and Logistic Regression (LR) classifiers from scratch, extended by non-linear kernel methods, on the standardized Wine Quality dataset. The study included comprehensive data preprocessing and rigorous hyperparameter tuning via 5-fold cross-validation, demonstrating a deep understanding of the models' underlying mechanisms. The key findings of this study are:

1. **Optimal Model Discrepancy (CV vs. Test):** The model that achieved the highest cross-validation score was the RBF-Kernel Logistic Regression (CV Accuracy: 0.7551). However, the final test set evaluation showed that the Linear-Kernel SVM delivered the highest overall performance across error metrics, achieving the best F1-Score of 0.8053 and the highest Recall (0.8290). This outcome confirms the SVM's robustness in generalization and its capability to maximize the margin on the final test distribution.
2. **Kernel Heterogeneity:** The optimal kernel selection was not universal. While the **\*\*RBF kernel\*\*** was crucial for the Logistic Regression model to achieve its peak performance, the **\*\*Linear kernel\*\*** yielded the highest scores for the SVM classifier. This suggests that the SVM's inherent mechanism was sufficient to find an effective separating hyperplane even in the original feature space, without relying on the non-linear RBF mapping, which often caused instability.

3. **Model Robustness and Convergence:** Both implemented models showed robust training dynamics, with rapid loss convergence and a minimal difference between training (Accuracy  $\approx 0.75 - 0.76$ ) and test accuracy (Accuracy  $\approx 0.73 - 0.74$ ). This confirms the successful control of overfitting through regularization ( $\lambda$ ) and effective feature standardization.