

# Protein Classification by Isoelectric Point Machine Learning Approach

Justyna Kowalska

Some plots are at the end of the file  
(All are in results folder)

## 1 Introduction

**Project Goal:** Classification of proteins and peptides into acidic ( $\text{pI} < 5$ ) and non-acidic ( $\text{pI} \geq 5$ ) using machine learning methods.

**Isoelectric Point (pI)** is the pH value at which a molecule carries no net electrical charge or is electrically neutral on average.

**Task:** Building ML models for binary classification of proteins/peptides based on their amino acid sequences.

## 2 Data Preparation

### 2.1 Data Sources

- **IPC peptide:** 16,882 items
- **IPC protein:** 2,324 items
- Source: <http://isoelectric.org/datasets.html>

---

Number of proteins with $\text{pI} < 5.0$	674
Number of peptides with $\text{pI} < 5.0$	5757
Number of proteins with $\text{pI} > 10.0$	19
Number of peptides with $\text{pI} > 10.0$	386

---

### 2.2 Preprocessing

- Labeling:  $\text{pI} < 5 \rightarrow \text{acidic (1)}$ ,  $\text{pI} \geq 5 \rightarrow \text{non-acidic (0)}$
- Split: 60% training / 20% testing / 20% validation

<b>Dataset</b>	<b>Total Items</b>	<b>Proteins</b>	<b>Peptides</b>	<b>Acidic</b>
Training	1001	347 (34.67%)	654 (65.33%)	663 (66.23%)
Testing	334	119 (35.63%)	215 (64.37%)	227 (67.96%)
Validation	334	103 (30.84%)	231 (69.16%)	210 (62.87%)

Table 1: Dataset distribution across training, testing, and validation sets

## 3 Features

### 3.1 Selected Features

Feature Name	Type	Description
uid	String	MD5 hash of sequence
data_type	String	Source type (protein/peptide)
pI	Float	Isoelectric point
length	Integer	Sequence length
molecular weight	Float	Molecular weight (Da)
label	Integer	Classification label (0/1)
aa_A, aa_C, ..., aa_Y	Integer	Count of each amino acid
hydro total	Float	Total hydrophobicity (Kyte-Doolittle)
hydro mean	Float	Mean hydrophobicity per residue
charge Sillero	Float	Net charge at pH 7 (Sillero scale)
charge EMBOSS	Float	Net charge at pH 7 (EMBOSS scale)
charge DTASelect	Float	Net charge at pH 7 (DTASelect scale)
charge Solomon	Float	Net charge at pH 7 (Solomon scale)
charge Rodwell	Float	Net charge at pH 7 (Rodwell scale)

Table 2: Feature descriptions

## 4 Machine Learning Models

### 4.1 Decision Tree

#### 4.1.1 Hyperparameter Optimization

A comprehensive grid search was performed with 5-fold cross-validation across 576 parameter combinations, totaling 2,880 model fits.

##### Best Parameters:

Parameter	Value
criterion	gini
max_depth	3
max_features	sqrt
min_samples_leaf	1
min_samples_split	2

Table 3: Optimal Decision Tree hyperparameters

Method	Best Parameter	CV Score	Validation Score
Depth Optimization	Depth: 1	0.997	0.997
Feature Selection	Features:1	0.997	0.997
Grid Search	Multiple params	0.997	1.000

Table 4: Decision Tree optimization methods comparison

#### 4.1.2 Performance Results

#### 4.1.3 Feature Importance Analysis

The Decision Tree model identified the most discriminative features for acidic protein classification:

Rank	Feature	Importance
1	charge Rodwell	0.9898
2	aa_R	0.0087
3	aa_N	0.0015
4-15	Other features	0.0000

Table 5: Top features by importance in Decision Tree model

The model demonstrates that **charge Rodwell** is by far the most important feature, accounting for 98.98% of the decision-making process, followed by arginine content (aa R) and asparagine content (aa N) with minimal contributions.

#### 4.1.4 Final Model Performance

##### Test Set Results:

- **Test Accuracy:** 99.4%
- **Cross-validation Score:** 99.7%
- **Validation Score:** 100.0%

Class	Precision	Recall	F1-score
Non-acidic	0.98	1.00	0.99
Acidic	1.00	0.99	1.00
<b>Weighted Avg</b>	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>

Table 6: Decision Tree classification report on test set

The Decision Tree model achieved exceptionally high performance metrics, what seems unusual to me.

## 4.2 Random Forest

The Random Forest model was optimized through systematic testing of multiple parameters, including number of estimators and comprehensive grid search across 648 parameter combinations (3,240 total fits with 5-fold CV).

### Number of Trees Optimization:

Trees	Train Acc	Val Acc	CV Score
10	1.000	1.000	0.997
25	1.000	1.000	0.997
50	1.000	1.000	0.997
100	1.000	1.000	0.997
200	1.000	1.000	0.997
500	1.000	1.000	0.997

Table 7: Random Forest n estimators optimization

### Best Parameters:

Parameter	Value
n_estimators	100
max_depth	5
max_features	sqrt
min_samples split	2
min_samples leaf	1
bootstrap	True

Table 8: Optimal Random Forest hyperparameters

### 4.2.1 Feature Importance Analysis

Unlike Decision Tree's single dominant feature, Random Forest shows more distributed importance:

Rank	Feature	Importance
1	charge Rodwell	0.226
2	charge Solomon	0.201
3	charge EMBOSS	0.183
4	charge Sillero	0.158
5	charge DTASelect	0.103
6	aa_E	0.036
7	length	0.026
8	aa_D	0.018

Table 9: Top Random Forest feature importance

#### 4.2.2 Ensemble Analysis

##### Tree Diversity Statistics:

- Average tree depth:  $2.60 \pm 1.02$
- Average nodes per tree:  $6 \pm 2$
- Feature importance correlation:  $0.069$
- Best individual tree accuracy: 100.0%
- Worst individual tree accuracy: 94.9%

#### 4.2.3 Final Performance

##### Test Set Results:

- **Test Accuracy:** 100.0%
- **Cross-validation Score:** 99.7%
- **Validation Score:** 100.0%

Class	Precision	Recall	F1-score
Non-acidic	1.00	1.00	1.00
Acidic	1.00	1.00	1.00
<b>Weighted Avg</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>

Table 10: Random Forest classification report

Random Forest achieved perfect test accuracy (100%). Suspiciously perfect performance continues to suggest potential data

#### 4.3 K-Nearest Neighbors (KNN)

The KNN model was analyzed through systematic testing of the number of neighbors (k), distance metrics, and feature selection methods. Data was standardized for optimal KNN performance.

Metric	CV Score
Euclidean	0.993
Manhattan	0.995
Chebyshev	0.957
Minkowski	0.995

Table 11: KNN distance metric comparison

Method	Optimal Features	CV Score
SelectKBest	12 features	1.000
RFE	8 features	0.999
Baseline (all)	30 features	0.995

Table 12: Feature selection comparison for KNN

#### 4.3.1 Best parameters

- algorithm=auto,
- metric=manhattan,
- n neighbors=7,
- weights=distance
- \*\*Cross-validation Score:\*\* 0.994
- \*\*Validation Accuracy:\*\* 0.997

#### 4.3.2 Performance Summary

Method	CV Score	Validation Score
K Optimization (k=1)	0.993	1.000
Distance Metrics (Manhattan)	0.995	-
Feature Selection (12 features)	1.000	-
Grid Search (Final)	0.994	0.997

Table 13: KNN optimization methods summary

#### 4.3.3 Final Model Evaluation

Class	Precision	Recall	F1-score
Non-acidic	1.00	0.99	1.00
Acidic	1.00	1.00	1.00
<b>Weighted Avg</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>

Table 14: KNN classification report on test set

Continues the pattern of suspiciously high performance across all models.

### 4.4 Support Vector Machine (SVM)

The SVM model was optimized through systematic testing of regularization parameter (C), gamma parameter and kernel types. Data was standardized for optimal SVM performance.

#### 4.4.1 Best parameters

- C=10
- gamma='scale'
- kernel='rbf'
- \*\*Cross-validation Score:\*\* 1.000
- \*\*Validation Accuracy:\*\* 1.000

#### 4.4.2 Feature Importance Analysis (Linear SVM)

Rank	Feature	Importance
1	charge DTASelect	2.027
2	charge EMBOSS	2.024
3	charge Sillero	2.022
4	charge Rodwell	1.993
5	charge Solomon	1.992
6	aa_H	0.796
7	aa_K	0.557
8	aa_C	0.505

Table 15: Top SVM feature importance (Linear kernel)

#### 4.4.3 Performance Summary

Method	CV Score	Validation Score
C Optimization (C=10)	1.000	-
Gamma Optimization (0.01)	0.999	-
Kernel Comparison (RBF)	1.000	-
Grid Search (Final)	1.000	1.000

Table 16: SVM optimization methods summary

#### 4.4.4 Final Model Evaluation

Class	Precision	Recall	F1-score
Non-acidic	1.00	1.00	1.00
Acidic	1.00	1.00	1.00
<b>Weighted Avg</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>

Table 17: SVM classification report on test set

Perfect performance continues the concerning pattern across all algorithms

## 4.5 Deep Learning Models

### 4.5.1 Model Architectures and Training Results

Ten different deep learning architectures were trained and evaluated on the dataset using CPU processing. All models achieved perfect accuracy (1.0000) on the training data, demonstrating the dataset's high separability.

Architecture	Layers	Parameters	Accuracy	Activation	BatchNorm	Dropout
Deep	5	51,905	1.0000	relu	True	0.3
Very_Deep	6	191,937	1.0000	relu	True	0.4
Wide	3	148,481	1.0000	relu	True	0.3
ELU Deep	5	51,905	1.0000	elu	True	0.3
Shallow	3	4,225	1.0000	relu	True	0.2
Medium	4	14,657	1.0000	relu	True	0.3
GELU Medium	4	14,657	1.0000	gelu	True	0.2
No_BatchNorm	4	14,209	1.0000	relu	False	0.4
High_Dropout	4	49,793	1.0000	relu	True	0.5
Ultra_Deep	8	139,265	1.0000	relu	True	0.3

Table 18: Comparison of deep learning architectures trained on the dataset

### 4.5.2 Architecture Analysis

#### Model Complexity Range:

- **Simplest:** Shallow (3 layers, 4,225 parameters)
- **Most Complex:** Very\_Deep (6 layers, 191,937 parameters)
- **Deepest:** Ultra\_Deep (8 layers, 139,265 parameters)

#### Activation Function Testing:

- **ReLU:** 7 models (standard choice)
- **ELU:** 1 model (ELU Deep)
- **GELU:** 1 model (GELU Medium)

#### Regularization Techniques:

- **Batch Normalization:** Used in 9/10 models
- **Dropout rates:** Ranging from 0.2 to 0.5
- **No BatchNorm model:** Tested impact of batch normalization

#### 4.5.3 Best Model Selection

The **Deep** architecture was selected as the best model based on its balanced complexity and performance:

##### Architecture Details:

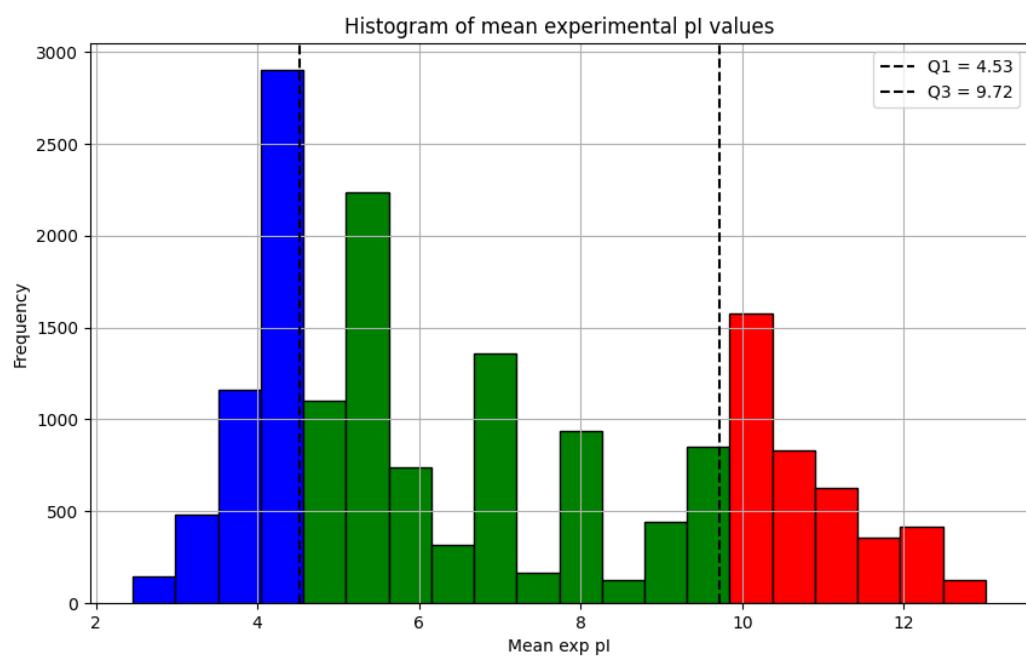
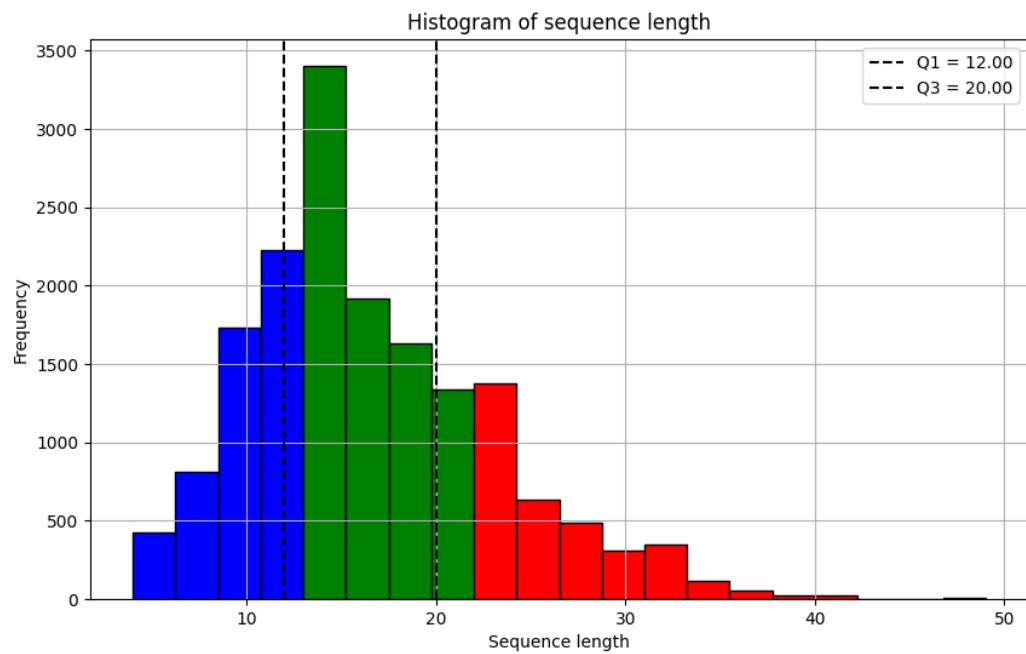
- **Configuration:** {'hidden dims': [256, 128, 64, 32], 'activation': 'relu', 'use batch norm': True, 'dropout rate': 0.3}
- **Parameters:** 51,905
- **Layers:** 5 (input + 4 hidden + output)
- **Final Accuracy:** 1.0000

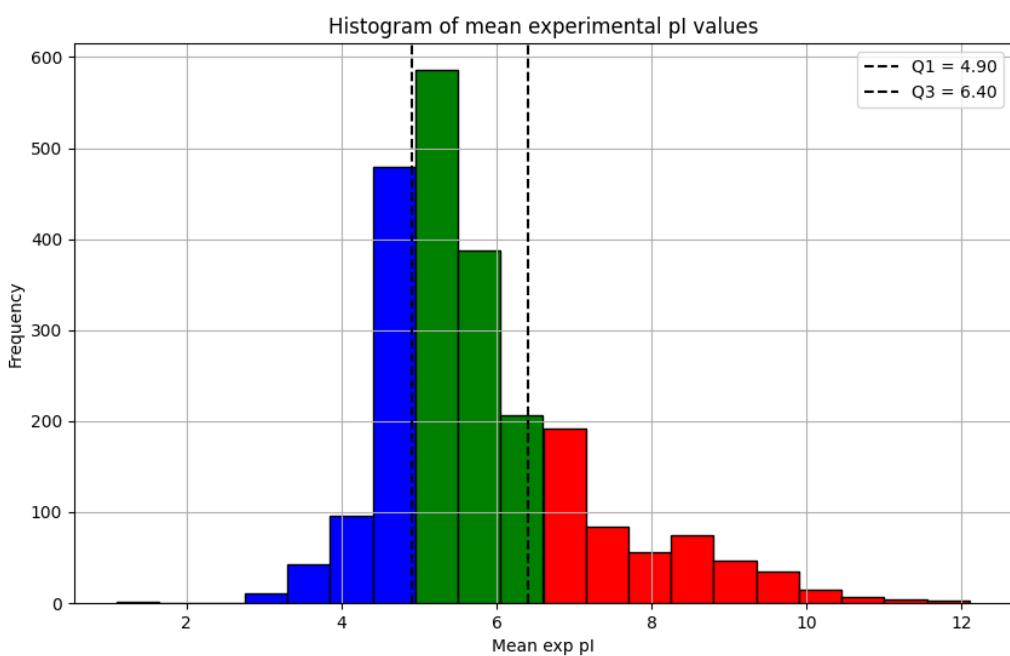
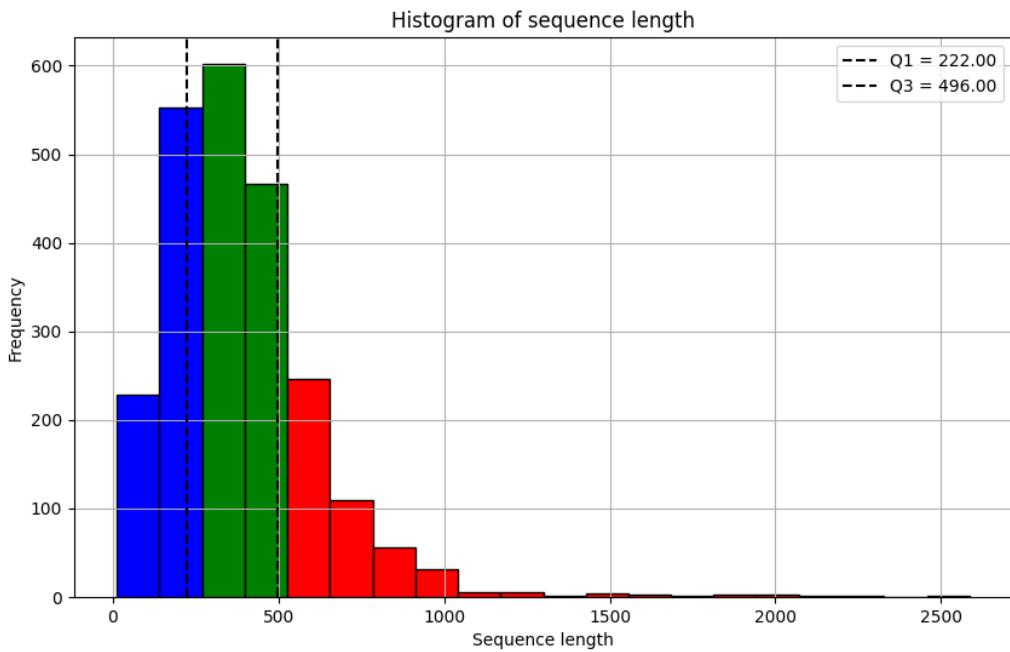
##### Perfect Performance Concerns:

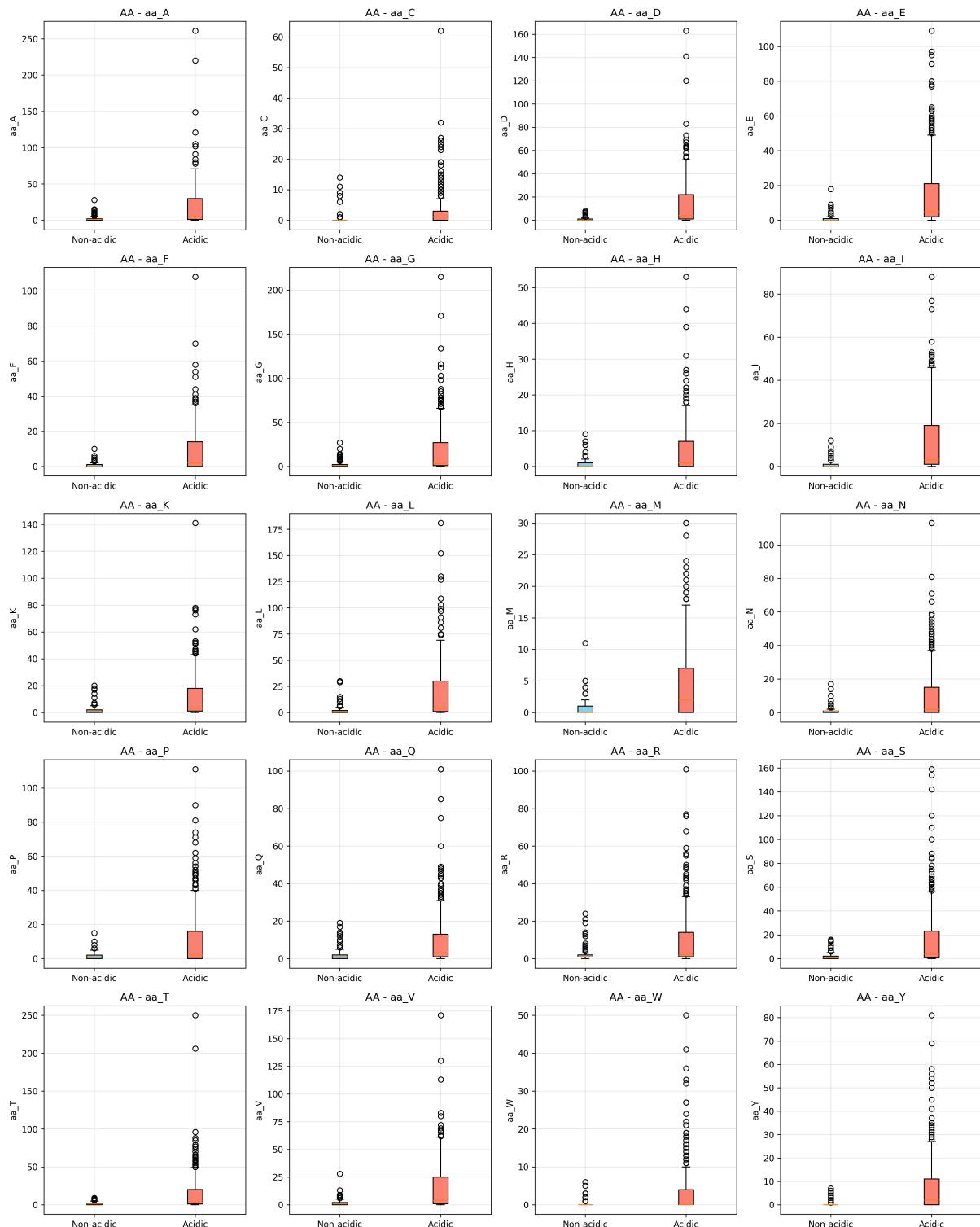
- All 10 architectures achieved identical perfect accuracy (1.0000)
- Performance independence from model complexity suggests dataset oversimplification
- Even the simplest Shallow model (4,225 parameters) achieved perfect results
- No BatchNorm model performed equally well, indicating strong feature separability

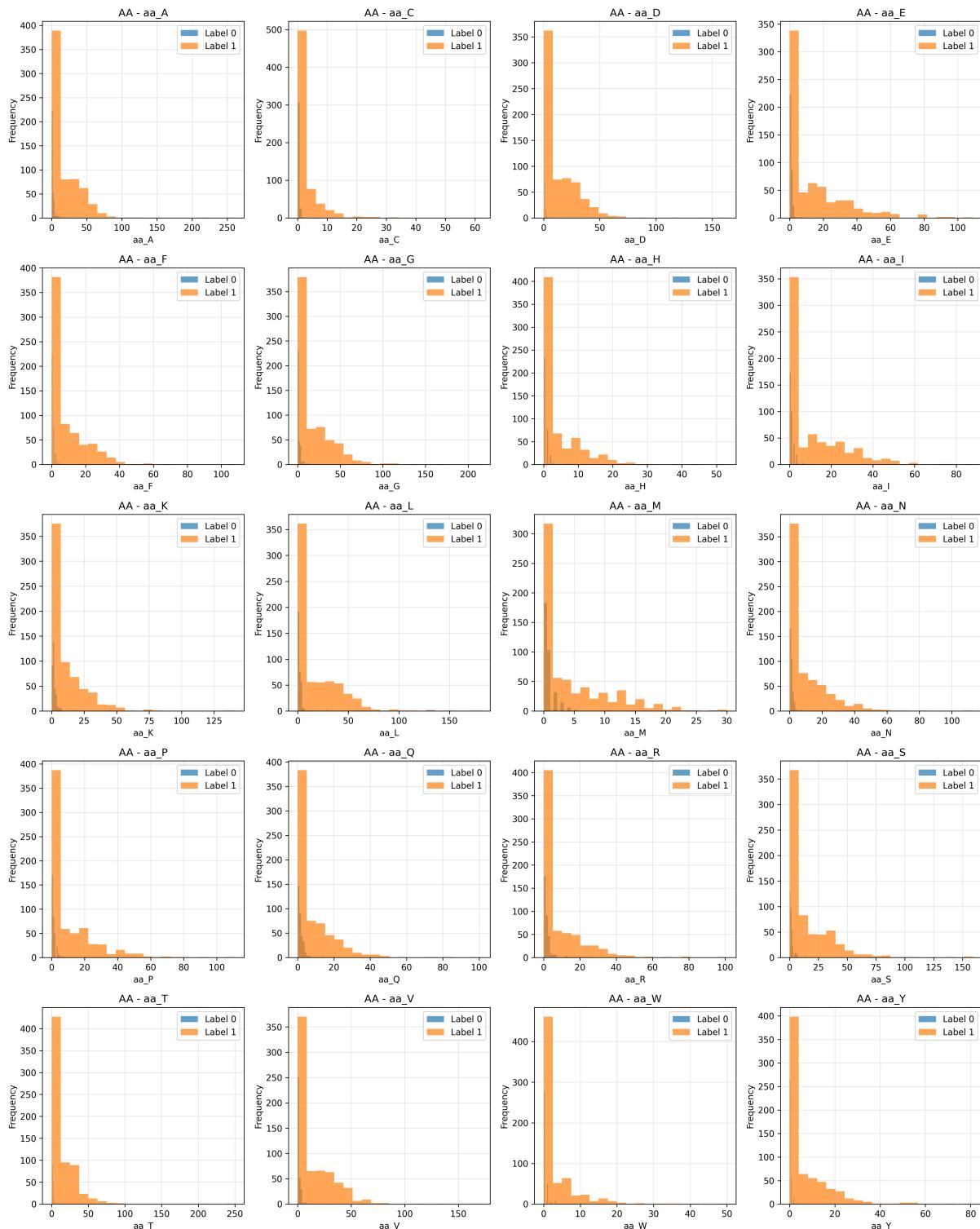
##### Model Robustness:

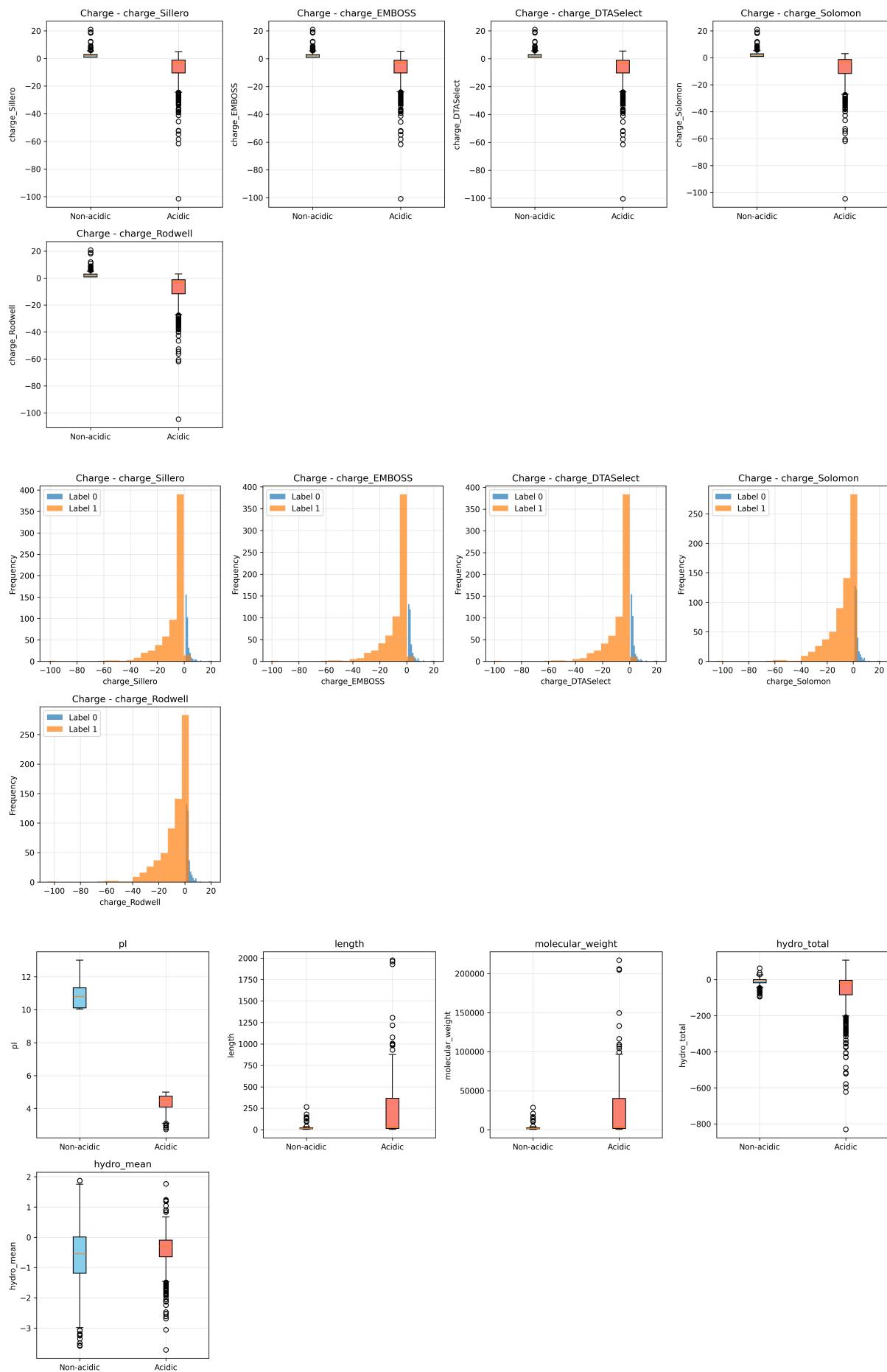
- Different activation functions (ReLU, ELU, GELU) showed no performance difference
- Varying dropout rates (0.2-0.5) had no impact on final accuracy
- Architecture depth (3-8 layers) did not affect classification performance

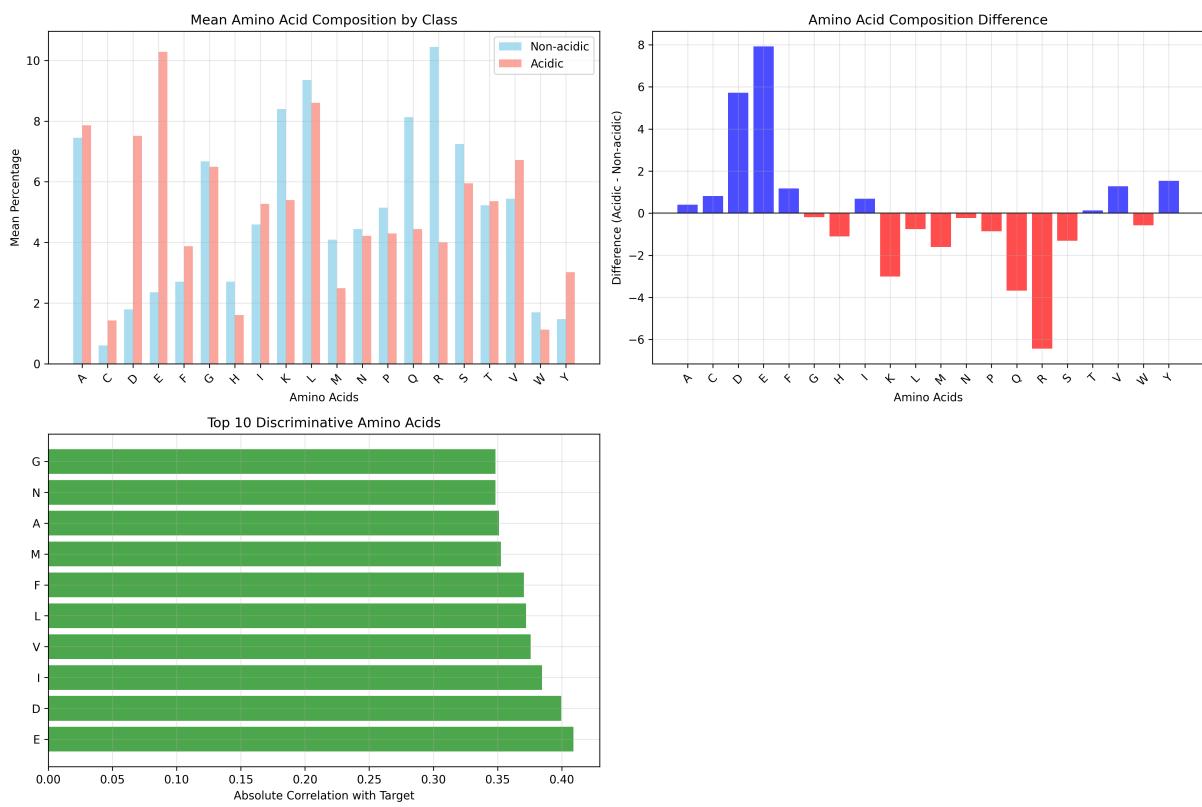
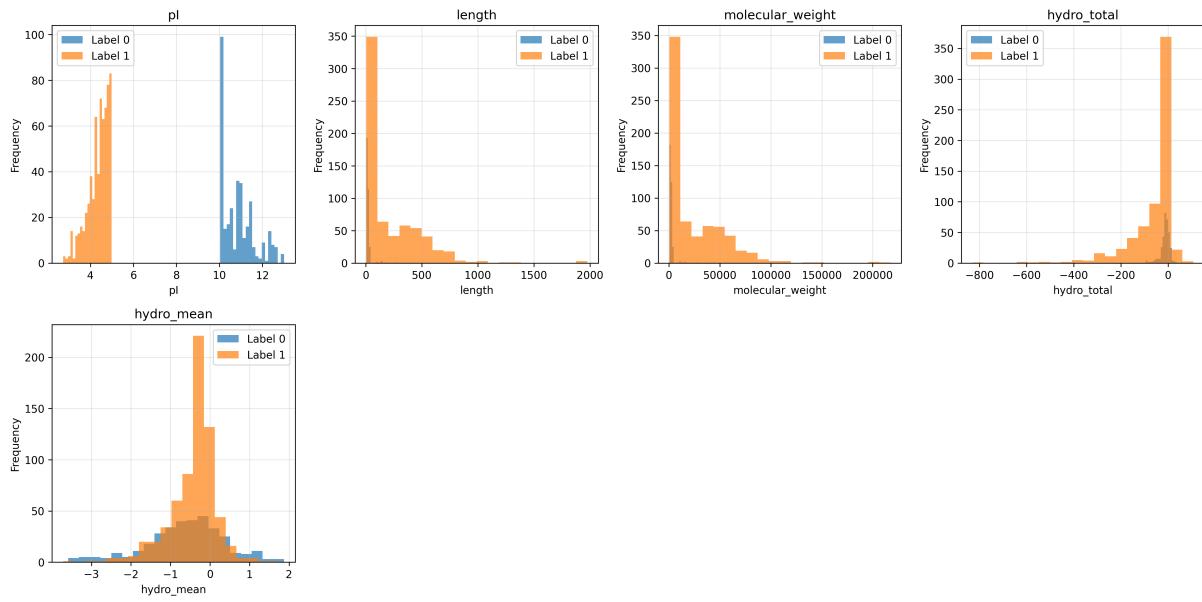


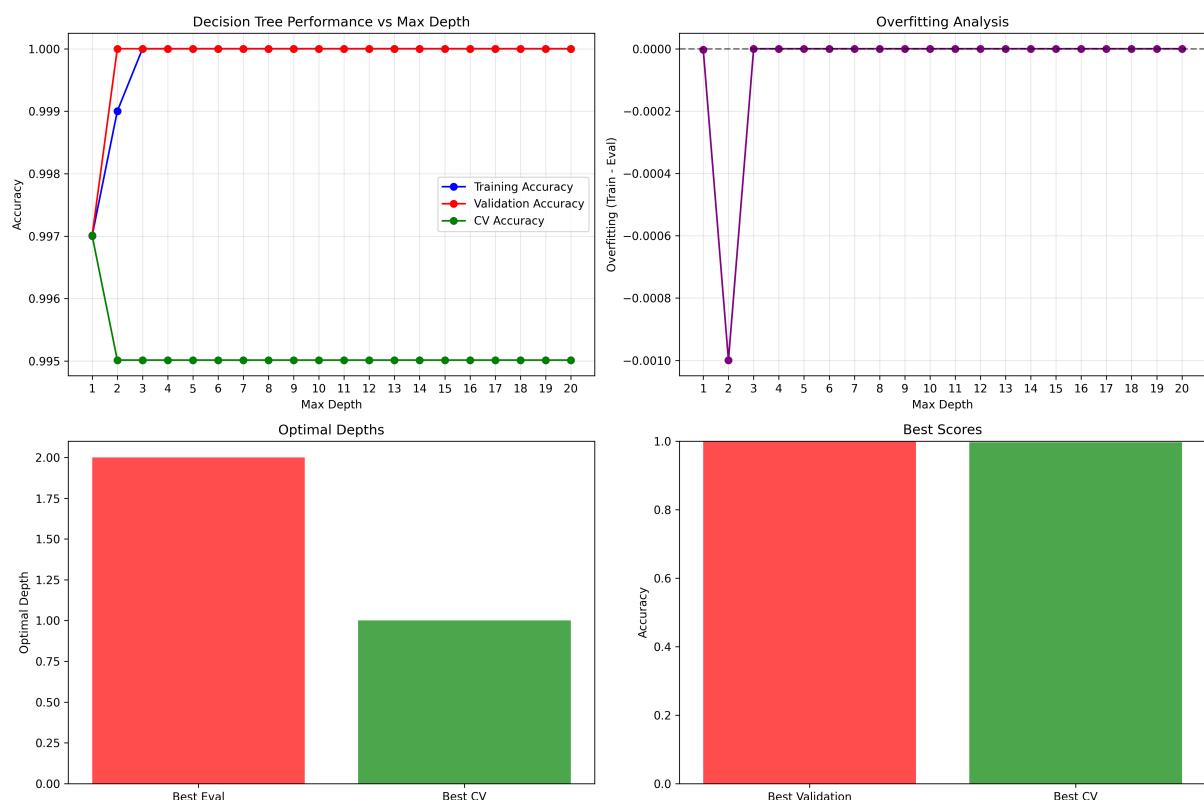
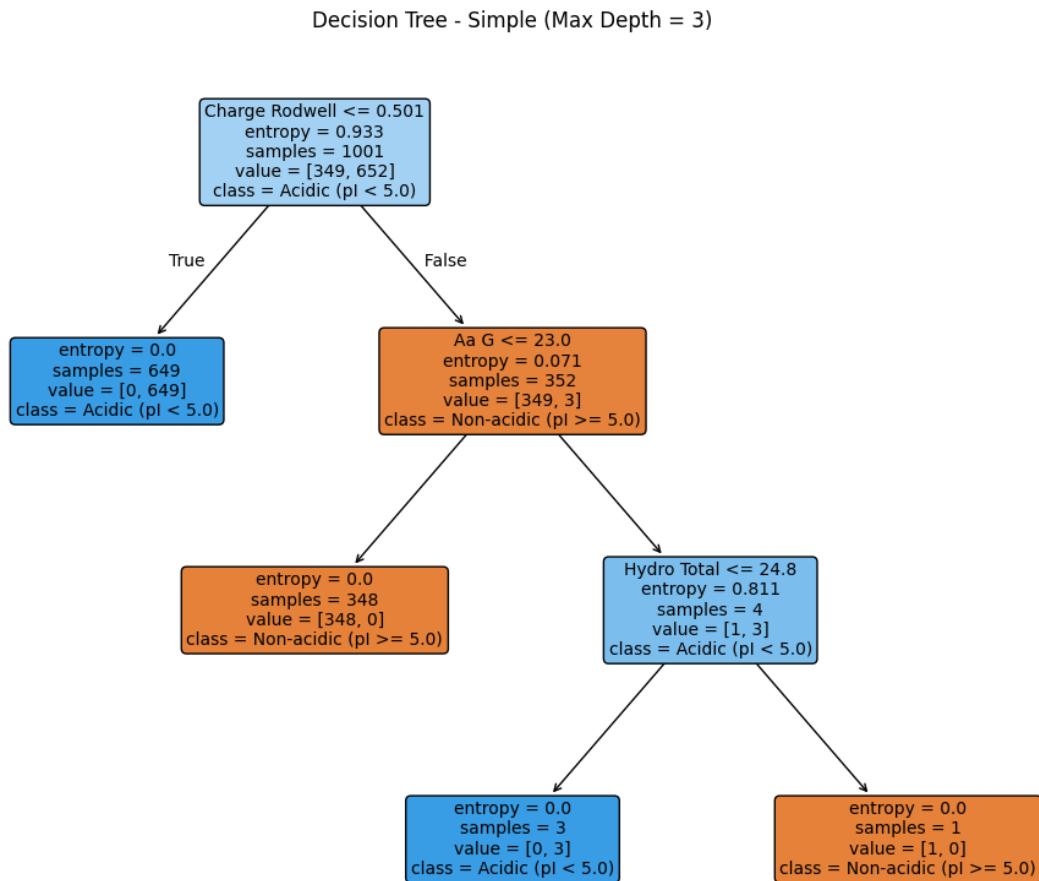




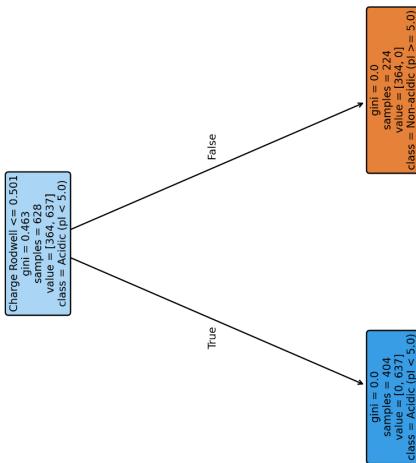




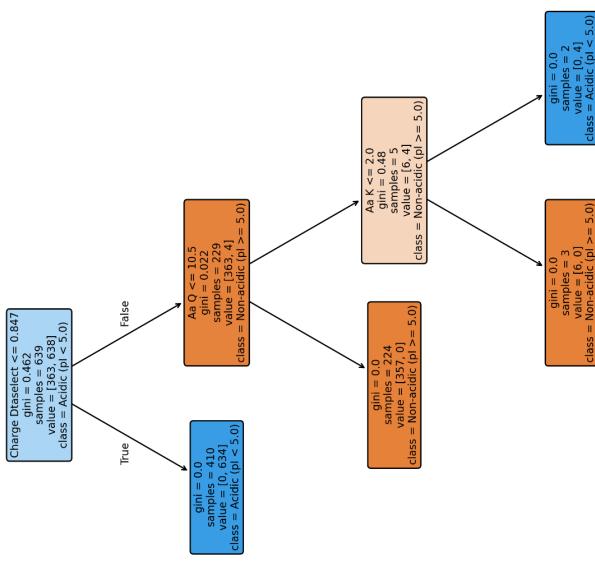




Tree 3 from Random Forest



Tree 2 from Random Forest



Tree 1 from Random Forest

