

---

# Comparing Supervised Learning Algorithms

---

Cecilia Truong  
COGS 118A  
13 December 2024

## Abstract

This work investigates the performance of three supervised learning algorithms (SVM with RBF kernel, KNN, and Decision Tree) on three datasets: Iris, Wine, and Breast Cancer (Wisconsin). The datasets were preprocessed and partitioned into three training and testing splits (80/20, 50/50, and 20/80) to evaluate the classifiers' accuracy under varying conditions. The results demonstrate that SVM consistently performed with the highest accuracy across all datasets. Decision Tree performed closely, often nearing SVM's performance, while KNN lagged behind, especially on datasets with greater complexity. These findings align with the trends observed in the empirical analysis by Caruana and Niculescu-Mizil, but refinements in the methodology to address overfitting could improve alignment with prior research.

## 1 Introduction

This paper investigates the accuracy performances of three supervised learning classifiers (**SVM with RBF kernel, KNN, and Decision Tree**) across three well-known datasets sourced from UCI's repository: **Iris, Wine, and Breast Cancer Wisconsin (Diagnostic)**. This study is inspired by Caruana and Niculescu-Mizil's work on the empirical analysis of supervised learning algorithms, which evaluated a wide array of models on multiple datasets. While their work was extensive, this report applies the core ideas on a smaller scale, focusing on three classifiers and three datasets. By doing so, it aims to provide an in-depth application of the concepts learned in COGS 118A and examine the comparative performance of these classifiers. The goal of this paper is to align with the findings of Caruana and Niculescu-Mizil while offering a concise and focused investigation. It seeks to analyze the accuracy of the classifiers under consistent experimental conditions, providing information on their relative strengths and weaknesses.

## 2 Problem & Data Description

The primary goal of this paper is to evaluate and compare the performance of three supervised learning classifiers on selected datasets. This study aims to align with the findings of Caruana and Niculescu-Mizil's empirical analysis of supervised learning algorithms, with some expected deviations due to differences in dataset selection and experimental setup. While the datasets are sourced using scikit-learn's library, they correspond directly to the versions available in the UCI Machine Learning Repository.

### 2.1 Iris Dataset

The Iris dataset is a simple, tabular dataset frequently used throughout COGS 118A. It contains 150 instances and 4 continuous features: sepal length, sepal width, petal length, and petal width. Each instance represents a plant, and the target variable classifies it into one of three species: Iris Setosa, Iris Versicolour, or Iris Virginica. Due to its simplicity and well-separated classes, this dataset serves as a baseline for evaluating the classifiers.

## 36 2.2 Wine Dataset

37 The Wine dataset includes chemical analysis results of wines grown in the same region of Italy but  
38 derived from three different cultivars. It comprises 178 instances and 13 continuous features, such as  
39 alcohol content, flavanoids, and color intensity. The target variable categorizes each instance into one  
40 of three wine classes. This dataset introduces higher dimensionality and slightly imbalanced classes,  
41 making it more challenging than the Iris dataset.

## 42 2.3 Breast Cancer (Wisconsin) Dataset

43 The Breast Cancer dataset is a multivariate dataset used to classify breast masses as malignant  
44 or benign based on characteristics of cell nuclei extracted from digitized images. It contains 569  
45 instances and 30 features. These 30 features are derived from 10 unique measurements (radius,  
46 texture, perimeter, area, compactness, etc.), with each measurement represented as the mean value,  
47 standard error, and worst (maximum) value. The target variable represents one of two diagnoses:  
48 malignant or benign. This structure provides a comprehensive view of the cell nuclei's characteristics  
49 and is the most complex out of the three datasets.

# 50 3 Methodology

51 Before implementing the classifiers, the datasets were preprocessed to simplify their class structures.  
52 For the Wine and Breast Cancer datasets, which originally had more than two classes, the classes  
53 were merged into two groups. Afterward, each dataset was divided into three distinct training and  
54 testing splits: 80/20, 50/50, and 20/80.

## 55 3.1 Classifiers

56 **SVM with RBF Kernel:** The Support Vector Machine classifier uses a radial basis function kernel  
57 to model non-linear relationships in the data. The hyperparameter gamma values ranged from  $10^{-6}$   
58 to  $10^{-2}$ , and C values ranged from 1 to 10,000. For each combination of C and gamma, 3-fold  
59 cross-validation was used to calculate average validation accuracy and training error. The combination  
60 that yielded the highest validation accuracy was selected as the optimal configuration.

61 **KNN:** The KNN classifier predicts the class of an instance based on the majority vote/average of  
62 its nearest neighbors. The number of neighbors,  $k$ , is the primary hyperparameter and was optimized  
63 by evaluating three values:  $k=1,2,3$ . For each  $k$ , 3-fold cross-validation was used to calculate the  
64 average validation accuracy and training error. The optimal  $k$  was selected based on the lowest  
65 average validation error.

66 **Decision Tree:** The Decision Tree classifier's key hyperparameter, max\_depth, was optimized  
67 through a grid search conducted over the values [1, 2, 3, 4, 5] using 3-fold cross-validation. The  
68 training and validation accuracies for each depth were recorded, and the most optimal max\_depth  
69 was selected based on the highest validation accuracy.

# 70 4 Experiments

71 **Iris Dataset:** On the Iris dataset, SVM achieved the highest accuracy of 0.9867, followed by  
72 Decision Tree at 0.9667 and KNN at 0.9467. The dataset's simplicity and well-separated classes  
73 made it favorable for SVM, which excels at creating smooth decision boundaries. Decision Tree also  
74 performed well, capturing the dataset's structure effectively. KNN performed slightly worse, but  
75 nonetheless on par with the rest of the classifiers.

76 **Breast Cancer(Wisconsin) Dataset:** Both SVM and Decision Tree achieved the highest accuracy  
77 of 0.9561 on the Breast Cancer dataset, significantly outperforming KNN, which reached 0.8421.  
78 This dataset's high dimensionality highlighted SVM's ability to handle complex feature interactions  
79 and Decision Tree's capacity to model patterns effectively. KNN struggled in performance most  
80 likely due to its incompatibility with complexity, as its distance-based approach becomes less reliable  
81 with a large number of features.

82 **Wine Dataset:** On the Wine dataset, SVM again outperformed the other classifiers with an accuracy  
83 of 0.9775, followed closely by Decision Tree at 0.9722. KNN had an accuracy of 0.8056, falling  
84 behind the other two classifiers. Like in the Breast Cancer results, the greater feature complexity and  
85 slight class imbalance of the Wine dataset demonstrated SVM's and Decision Tree's flexibility, while  
86 KNN's struggles were consistent with its sensitivity to these dataset characteristics.

87 The results of this investigation closely align with the findings of Caruana and Niculescu-Mizil's  
88 empirical analysis. Among the three supervised learning algorithms evaluated, SVM consistently  
89 had the highest performance in both this study and Caruana and Niculescu-Mizil's work. However,  
90 there were noteworthy deviations in the performances of KNN and Decision Tree compared to the  
91 original study. Specifically, Decision Tree significantly outperformed KNN on the more complex  
92 datasets, such as Wine and Breast Cancer, despite falling short of KNN's performance in Caruana  
93 and Niculescu-Mizil's empirical analysis.

## 94 **5 Conclusion**

95 This study compared the performance of SVM with RBF kernel, KNN, and Decision Tree across three  
96 datasets (Iris, Wine, and Breast Cancer Wisconsin). The results demonstrated that SVM consistently  
97 achieved the highest accuracy, with the Decision Tree performing nearly as well as SVM in most  
98 cases. KNN, while effective on simpler datasets like Iris, underperformed on complex datasets. The  
99 findings could potentially align more closely with Caruana and Niculescu-Mizil's work through  
100 improved strategies for addressing overfitting.

## 101 **References**

- 102 [1] Caruana, R. & Niculescu-Mizil, A. (2006) An empirical comparison of supervised learning algorithms. In  
103 Proceedings of the 23rd international conference on Machine learning (ICML '06). Association for Computing  
104 Machinery, New York, NY, USA, 161–168. <https://doi.org/10.1145/1143844.1143865>
- 105 [2] Dua, D., & Graff, C. (2019). UCI Machine Learning Repository. University of California, Irvine, School of  
106 Information and Computer Sciences. <https://archive.ics.uci.edu/ml>
- 107 [3] Fisher, R. (1936). Iris [Dataset]. UCI Machine Learning Repository. <https://doi.org/10.24432/C56C76>.
- 108 [4] Aeberhard, S. & Forina, M. (1992). Wine [Dataset]. UCI Machine Learning Repository.  
109 <https://doi.org/10.24432/C5PC7J>.
- 110 [5] Wolberg, W., Mangasarian, O., Street, N., & Street, W. (1993). Breast Cancer Wisconsin (Diagnostic)  
111 [Dataset]. UCI Machine Learning Repository. <https://doi.org/10.24432/C5DW2B>.