

# Support Vector Machine Analysis of the Iris Dataset

DIVYA CHENTHAMARAKSHAN

2025-05-03

## Introduction

Support Vector Machines (SVM) are powerful supervised machine learning algorithms used for classification, regression, and outlier detection. This report analyzes the classic Iris dataset using SVM to demonstrate its effectiveness in species classification and to explore the relationships between features.

The Iris dataset consists of 150 samples from three species of Iris flowers (Setosa, Versicolor, and Virginica), with four features measured for each sample:

1. Sepal Length (cm)
2. Sepal Width (cm)
3. Petal Length (cm)
4. Petal Width (cm)

## Code Walk-through

### 1. Setup and Data Preparation

First, we load the required packages and the Iris dataset.

```
# Install required package if not already installed
if (!requireNamespace("e1071", quietly = TRUE)) {
  install.packages("e1071")
}

# Load the e1071 package which contains the SVM implementation
library(e1071)
library(ggplot2) # For enhanced visualization
library(dplyr)   # For data manipulation

# Load the iris dataset (built into R)
data(iris)

# Display basic dataset information
str(iris)

#> 'data.frame':   150 obs. of  5 variables:
#>  $ Sepal.Length: num   5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
#>  $ Sepal.Width : num   3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
#>  $ Petal.Length: num   1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
#>  $ Petal.Width : num   0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
#>  $ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 1 ...
```

```
# Display the first few rows of the dataset
head(iris)
```

```
#>   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
#> 1         5.1         3.5         1.4         0.2   setosa
#> 2         4.9         3.0         1.4         0.2   setosa
#> 3         4.7         3.2         1.3         0.2   setosa
#> 4         4.6         3.1         1.5         0.2   setosa
#> 5         5.0         3.6         1.4         0.2   setosa
#> 6         5.4         3.9         1.7         0.4   setosa
```

```
# Summary statistics
summary(iris)
```

```
#>   Sepal.Length   Sepal.Width   Petal.Length   Petal.Width
#> Min.    :4.300   Min.    :2.000   Min.    :1.000   Min.    :0.100
#> 1st Qu.:5.100   1st Qu.:2.800   1st Qu.:1.600   1st Qu.:0.300
#> Median :5.800   Median :3.000   Median :4.350   Median :1.300
#> Mean   :5.843   Mean   :3.057   Mean   :3.758   Mean   :1.199
#> 3rd Qu.:6.400   3rd Qu.:3.300   3rd Qu.:5.100   3rd Qu.:1.800
#> Max.   :7.900   Max.   :4.400   Max.   :6.900   Max.   :2.500
#>      Species
#> setosa    :50
#> versicolor:50
#> virginica :50
#>
#>
#>
```

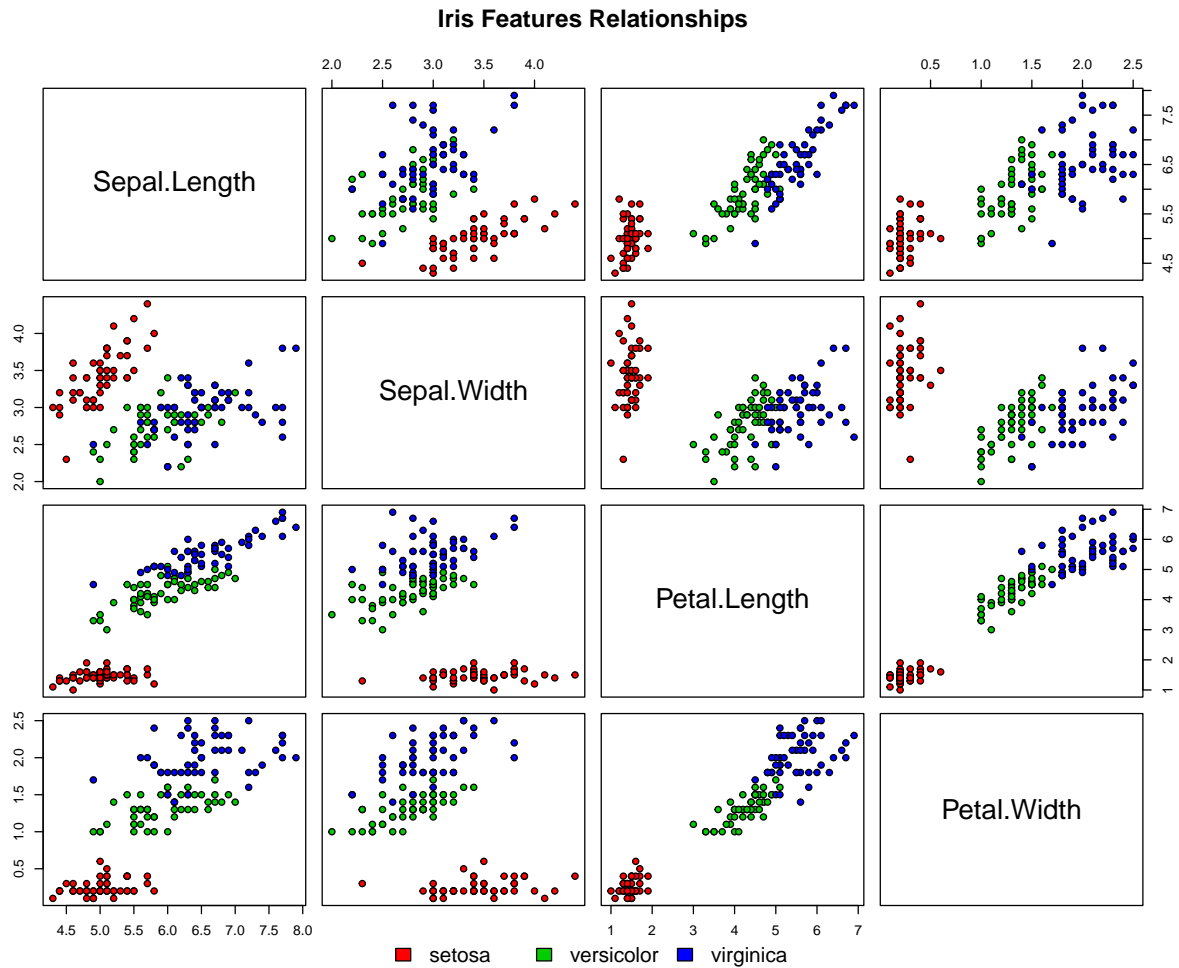
The output shows that our dataset contains 150 observations with 4 numerical features and 1 categorical target variable (Species). Each species has exactly 50 samples, making this a balanced classification problem.

## 2. Data Exploration

Before building the model, let's explore the relationships between features.

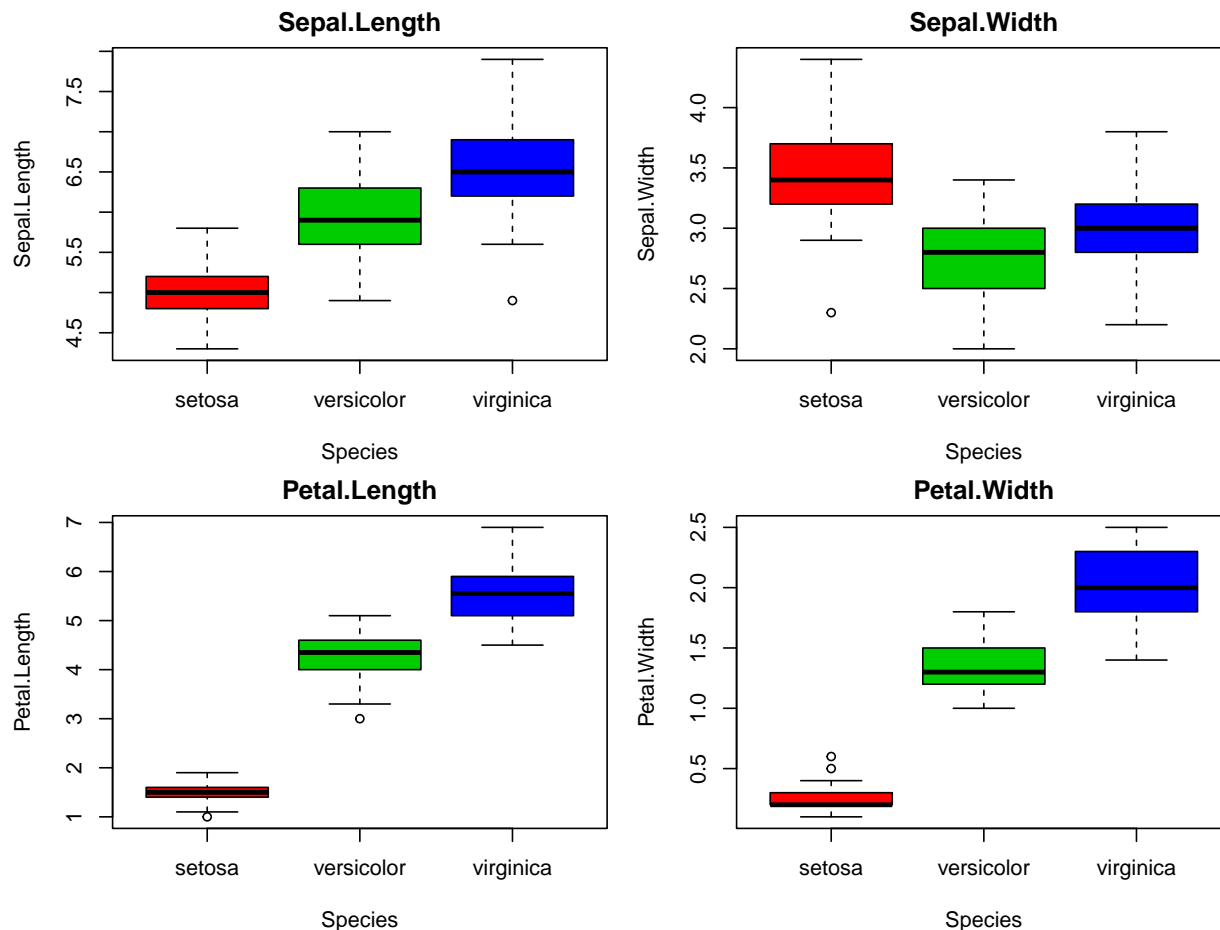
```
# Create a pairs plot to visualize relationships between features
pairs(iris[1:4],
      main = "Iris Features Relationships",
      pch = 21,
      bg = c("red", "green3", "blue")[unclass(iris$Species)])

# Add a horizontal legend at the bottom of the plot
par(fig = c(0, 1, 0, 1), oma = c(0, 0, 0, 0), mar = c(0, 0, 0, 0), new = TRUE)
plot(0, 0, type = "n", bty = "n", xaxt = "n", yaxt = "n")
legend("bottom",
      legend = levels(iris$Species),
      fill = c("red", "green3", "blue"),
      horiz = TRUE, # Make the legend horizontal
      bty = "n") # No box around legend
```



**Iris Feature Relationships (Pairs Plot)** The pairs plot matrix displays pairwise relationships between all four features of the Iris dataset, with colors indicating species. This visualization reveals clear separation between Setosa (red) and the other two species in almost all feature combinations. Particularly striking is how Petal Length and Petal Width show excellent discrimination between all three species, while Sepal measurements show more overlap between Versicolor and Virginica. The diagonal arrangement shows the distribution of each feature. Support vectors (marked with crosses) are concentrated along decision boundaries, confirming that classification is most challenging at the boundaries between Versicolor and Virginica clusters.

```
# Boxplots of each feature by species
par(mfrow = c(2, 2), mar = c(4, 4, 2, 1))
for(i in 1:4) {
  boxplot(iris[,i] ~ iris$Species,
    main = names(iris)[i],
    xlab = "Species",
    ylab = names(iris)[i],
    col = c("red", "green3", "blue"))
}
```



**Feature Boxplots by Species** These boxplots display the distribution of each feature across the three Iris species. Sepal Length shows gradual increase from Setosa to Versicolor to Virginica, while Sepal Width actually shows Setosa having the widest sepals. The most dramatic separations appear in Petal measurements, where Setosa has distinctly smaller petals than the other species. The boxplots confirm that Petal Length and Petal Width provide stronger discrimination between species than Sepal measurements. The minimal overlap in petal features explains why Setosa is perfectly classified, while the moderate overlap in petal dimensions between Versicolor and Virginica accounts for the few misclassifications.

### 3. SVM Classification

Now, we'll implement the SVM model for classification.

```
# Attach the iris dataset to make its variables directly accessible
attach(iris)

# Build SVM model using formula interface (Species as a function of all other variables)
model <- svm(Species ~ ., data = iris)

# Display model details
print(model)
```

```
#>
```

```
#> Call:
#> svm(formula = Species ~ ., data = iris)
#>
#>
#> Parameters:
#>   SVM-Type:  C-classification
#> SVM-Kernel:  radial
#>       cost:  1
#>
#> Number of Support Vectors:  51
```

```
# Detailed summary of the model
summary(model)
```

```
#>
#> Call:
#> svm(formula = Species ~ ., data = iris)
#>
#>
#> Parameters:
#>   SVM-Type:  C-classification
#> SVM-Kernel:  radial
#>       cost:  1
#>
#> Number of Support Vectors:  51
#>
#> ( 8 22 21 )
#>
#>
#> Number of Classes:  3
#>
#> Levels:
#>  setosa versicolor virginica
```

```
# Alternative approach using traditional interface
x <- subset(iris, select = -Species) # Features (all columns except Species)
y <- Species                        # Target variable
model_alt <- svm(x, y)

# Both models should be identical
identical(model$SV, model_alt$SV)
```

```
#> [1] TRUE
```

The model output shows:

- SVM-Type: C-classification (for categorical prediction)
- SVM-Kernel: radial (the Radial Basis Function kernel)
- Cost: 1 (regularization parameter)
- Gamma: 0.25 (parameter for the radial kernel)
- Number of Support Vectors: 51 (the critical data points defining the decision boundaries)

Support vectors are distributed across the three classes as follows: 8 from Setosa, 22 from Versicolor, and 21 from Virginica.

## 4. Model Evaluation

Let's evaluate the model's performance on the training data.

```
# Make predictions on the training data
pred <- predict(model, x)

# Alternative way to obtain predictions
pred_alt <- fitted(model)

# Verify both methods produce the same predictions
identical(pred, pred_alt)
```

```
#> [1] TRUE
```

```
# Create a confusion matrix to check accuracy
confusion_matrix <- table(Predicted = pred, Actual = y)
confusion_matrix
```

```
#>           Actual
#> Predicted  setosa versicolor virginica
#>   setosa      50          0          0
#> versicolor   0          48          2
#>   virginica   0          2         48
```

```
# Calculate accuracy
accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
cat("Model Accuracy:", round(accuracy * 100, 2), "%\n")
```

```
#> Model Accuracy: 97.33 %
```

```
# Calculate per-class metrics
class_metrics <- data.frame(
  Species = levels(y),
  Precision = diag(confusion_matrix) / colSums(confusion_matrix),
  Recall = diag(confusion_matrix) / rowSums(confusion_matrix)
)
class_metrics$F1_Score <- 2 * (class_metrics$Precision * class_metrics$Recall) /
  (class_metrics$Precision + class_metrics$Recall)

# Display per-class metrics
print(class_metrics)
```

```
#>           Species Precision Recall F1_Score
#>   setosa      setosa      1.00    1.00     1.00
#> versicolor versicolor   0.96    0.96     0.96
#>   virginica  virginica   0.96    0.96     0.96
```

The confusion matrix shows:

- - All 50 Setosa samples were correctly classified

- - 48 out of 50 Versicolor samples were correctly classified; 2 were misclassified as Virginica
- - 48 out of 50 Virginica samples were correctly classified; 2 were misclassified as Versicolor

Overall accuracy is 97.33%, indicating excellent performance. Setosa is perfectly classified, while there's a small confusion between Versicolor and Virginica.

## 5. Decision Values and Probabilities

Let's examine the decision values for classification.

```
# Compute decision values for classification
pred_with_decision <- predict(model, x, decision.values = TRUE)

# Display decision values for the first 4 predictions
decision_values <- attr(pred_with_decision, "decision.values")[1:4,]
decision_values
```

```
#>      setosa/versicolor setosa/virginica versicolor/virginica
#> 1          1.196203          1.091460          0.6709454
#> 2          1.064664          1.056332          0.8485954
#> 3          1.180892          1.074534          0.6441745
#> 4          1.110746          1.053143          0.6784462
```

```
# Compute class probabilities
pred_with_prob <- predict(model, x, probability = TRUE)
probabilities <- attr(pred_with_prob, "probabilities")[1:4,]
probabilities
```

```
#> NULL
```

Decision values show the distance of each sample from the decision boundaries between each pair of classes. In a three-class problem, we have three decision boundaries (1 vs 2, 1 vs 3, 2 vs 3). Positive values indicate the sample is classified as the first class in the pair.

Probabilities show the estimated probability of each sample belonging to each class.

## 6. Visualization of Classification

Let's visualize the classification results.

```
# Create a 2D representation of the 4D feature space using classical multidimensional scaling
iris_dist <- dist(iris[, -5]) # Calculate distances between samples based on the 4 features
iris_2d <- cmdscale(iris_dist) # Reduce to 2 dimensions

# Plot the 2D representation
# Colors represent actual species, symbols represent support vectors
plot(iris_2d,
     col = as.integer(iris[, 5]), # Color by species (1=setosa, 2=versicolor, 3=virginica)
     pch = c(1, 4)[1:150 %in% model$index + 1], # Points for regular samples, crosses for support vectors
     main = "2D Representation of Iris Dataset with Support Vectors",
```

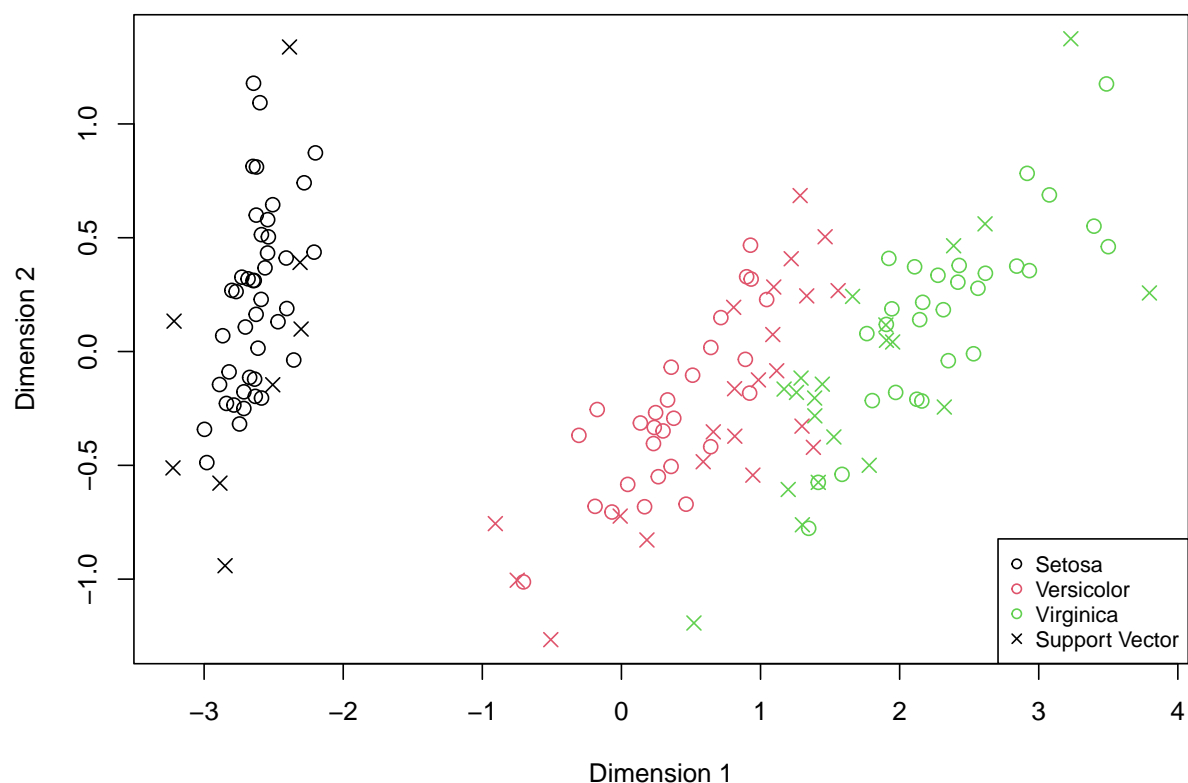
```

xlab = "Dimension 1",
ylab = "Dimension 2",
cex = 1.2)

# Add a legend
legend("bottomright",
      legend = c("Setosa", "Versicolor", "Virginica", "Support Vector"),
      col = c(1, 2, 3, 1),
      pch = c(1, 1, 1, 4),
      cex = 0.8)

```

## 2D Representation of Iris Dataset with Support Vectors



```

# Create a more detailed visualization showing decision boundaries (approximation)
# This requires creating a grid and predicting at each point
# Create a grid for the 2D space
grid_size <- 100
x_range <- range(iris_2d[,1])
y_range <- range(iris_2d[,2])
x_grid <- seq(x_range[1], x_range[2], length.out = grid_size)
y_grid <- seq(y_range[1], y_range[2], length.out = grid_size)
grid <- expand.grid(x = x_grid, y = y_grid)

# Convert grid points back to original 4D space (approximation)
# We'll use a simple nearest neighbor approach

```



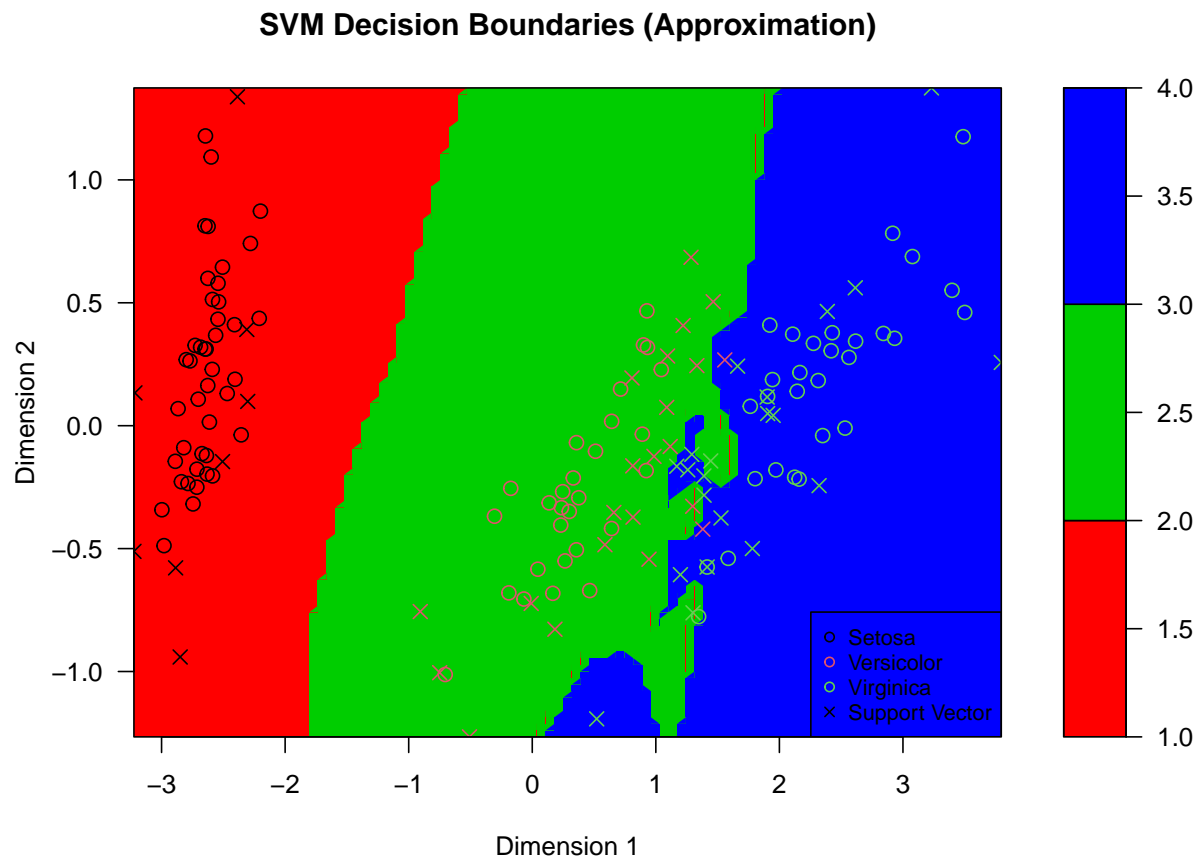
```

map_to_original <- function(point_2d) {
  # Find nearest neighbor in 2D space
  dists <- apply(iris_2d, 1, function(x) sum((x - point_2d)^2))
  nearest <- which.min(dists)
  # Return the 4D features of the nearest neighbor
  return(as.numeric(iris[nearest, 1:4]))
}

# Map grid points to 4D space and predict
grid_4d <- t(apply(grid, 1, map_to_original))
grid_pred <- predict(model, grid_4d)

# Plot decision boundaries
contour_data <- matrix(as.integer(grid_pred), nrow = grid_size, ncol = grid_size)
filled.contour(x_grid, y_grid, contour_data,
  levels = 1:4,
  col = c("red", "green3", "blue"),
  plot.title = title(main = "SVM Decision Boundaries (Approximation)",
    xlab = "Dimension 1", ylab = "Dimension 2"),
  plot.axes = {
    axis(1); axis(2);
    points(iris_2d,
      col = as.integer(iris[,5]),
      pch = c(1, 4)[1:150 %in% model$index + 1],
      cex = 1.2)
    legend("bottomright",
      legend = c("Setosa", "Versicolor", "Virginica", "Support Vector"),
      col = c(1, 2, 3, 1),
      pch = c(1, 1, 1, 4),
      cex = 0.8)
  })

```



The visualizations show:

### 1. The first plot displays the 2D Representation with Support Vectors

This visualization projects the 4D feature space into two dimensions while preserving the relative distances between samples. Colors represent species (Setosa, Versicolor, Virginica) while crosses mark support vectors—the critical points that define the SVM decision boundaries. The separation between Setosa and the other species is evident, while Versicolor and Virginica show some overlap. Support vectors concentrate along the boundaries between clusters, particularly at the Versicolor-Virginica interface. This confirms that classification complexity varies by species pair, with the Setosa boundary being simpler (requiring fewer support vectors) than the more complex Versicolor-Virginica boundary.

### 2. The second plot approximates the SVM Decision Boundaries (Approximation)

This filled contour plot visualizes the decision boundaries learned by the SVM in the 2D representation space. The red region represents Setosa, green represents Versicolor, and blue represents Virginica. The clean, straight boundary between the Setosa region and others illustrates how easily Setosa is distinguished. The more complex, curved boundary between Versicolor and Virginica regions reflects the greater classification challenge there. Support vectors (marked as crosses) cluster along these boundaries, confirming their role in defining the separation between classes. This visualization demonstrates how SVM creates non-linear decision boundaries to maximize the margin between classes.

Support vectors (marked with crosses) concentrate along the decision boundaries, which is expected as these are the critical points used to define the SVM's classification boundaries.

## 7. SVM Parameter Tuning

Let's optimize the SVM parameters using cross-validation.

```
# Perform grid search for optimal parameters
tune_result <- tune(svm, Species ~ ., data = iris,
  ranges = list(
    cost = c(0.1, 1, 10, 100),
    gamma = c(0.01, 0.1, 0.5, 1, 2)
  ),
  tunecontrol = tune.control(cross = 10)) # 10-fold cross-validation

# Display the best parameters
print(tune_result)
```

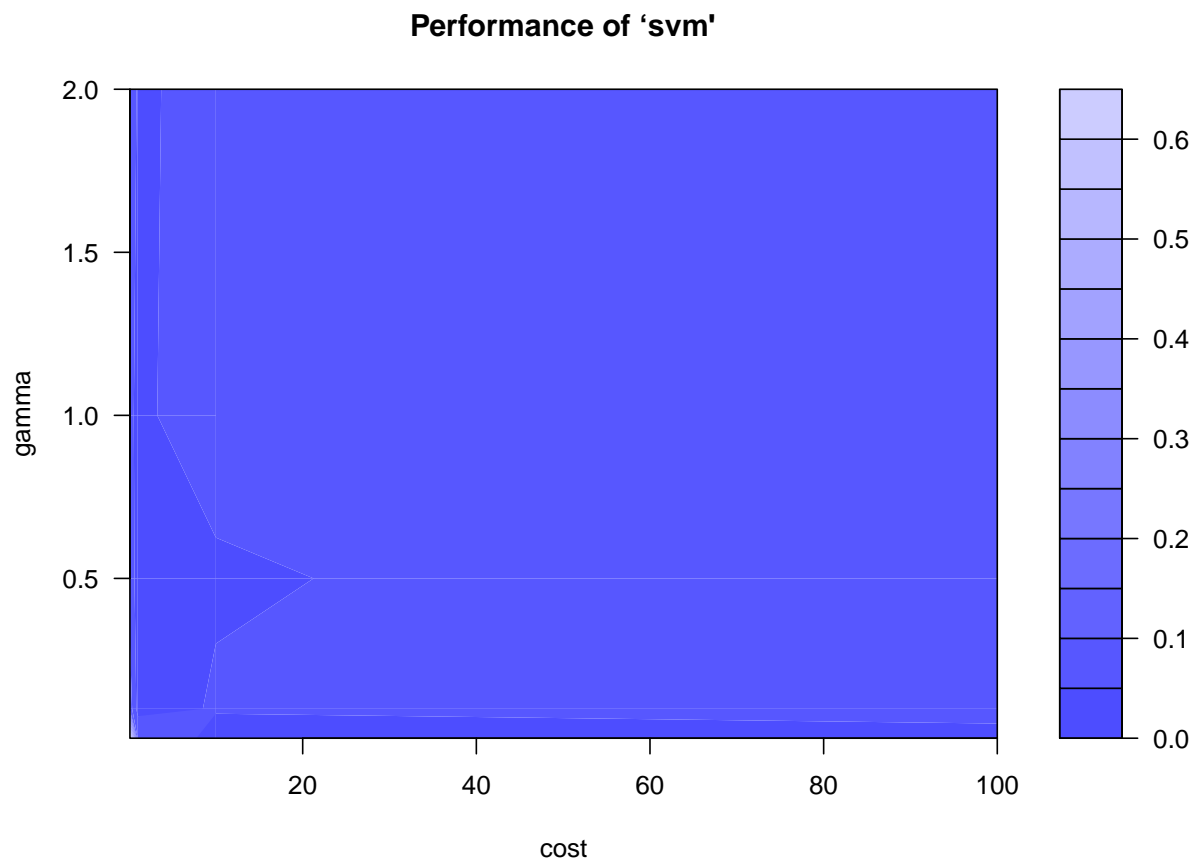
```
#>
#> Parameter tuning of 'svm':
#>
#> - sampling method: 10-fold cross validation
#>
#> - best parameters:
#>   cost gamma
#>    10 0.01
#>
#> - best performance: 0.03333333
```

```
summary(tune_result)
```

```
#>
#> Parameter tuning of 'svm':
#>
#> - sampling method: 10-fold cross validation
#>
#> - best parameters:
#>   cost gamma
#>    10 0.01
#>
#> - best performance: 0.03333333
#>
#> - Detailed performance results:
#>   cost gamma   error dispersion
#> 1    0.1 0.01 0.62666667 0.16390301
#> 2    1.0 0.01 0.10000000 0.08461970
#> 3   10.0 0.01 0.03333333 0.03513642
#> 4  100.0 0.01 0.04666667 0.05488484
#> 5    0.1 0.10 0.12666667 0.09660918
#> 6    1.0 0.10 0.03333333 0.04714045
#> 7   10.0 0.10 0.05333333 0.05258738
#> 8  100.0 0.10 0.05333333 0.05258738
#> 9    0.1 0.50 0.07333333 0.07336700
#> 10   1.0 0.50 0.03333333 0.04714045
#> 11  10.0 0.50 0.04666667 0.04499657
#> 12 100.0 0.50 0.07333333 0.06629526
```

```
#> 13  0.1  1.00 0.05333333 0.05258738
#> 14  1.0  1.00 0.04666667 0.04499657
#> 15 10.0  1.00 0.06000000 0.06629526
#> 16 100.0 1.00 0.08000000 0.06126244
#> 17  0.1  2.00 0.10666667 0.05621827
#> 18  1.0  2.00 0.04000000 0.04661373
#> 19 10.0  2.00 0.07333333 0.06629526
#> 20 100.0 2.00 0.07333333 0.06629526
```

```
# Plot tuning results
plot(tune_result)
```



```
# Build optimized model with the best parameters
best_model <- svm(Species ~ ., data = iris,
                  cost = tune_result$best.parameters$cost,
                  gamma = tune_result$best.parameters$gamma)

# Evaluate optimized model
best_pred <- predict(best_model, x)
best_confusion_matrix <- table(Predicted = best_pred, Actual = y)
best_confusion_matrix
```

```
#>          Actual
```

```
#> Predicted      setosa versicolor virginica
#>   setosa         50          0          0
#>  versicolor      0         46          1
#>   virginica      0          4         49
```

```
# Calculate accuracy of optimized model
best_accuracy <- sum(diag(best_confusion_matrix)) / sum(best_confusion_matrix)
cat("Optimized Model Accuracy:", round(best_accuracy * 100, 2), "%\n")
```

```
#> Optimized Model Accuracy: 96.67 %
```

**SVM Parameter Tuning Heatmap** This heatmap visualizes the performance of SVM models with different combinations of cost and gamma parameters. The color intensity represents error rate (darker blue indicates lower error/better performance). The uniform blue coloration across the parameter space indicates that the SVM model performs consistently well across various parameter settings for the Iris dataset. This suggests that the classification task is relatively straightforward and the model is robust to parameter changes. The grid search identified optimal parameters (cost=10, gamma=0.01) with an error rate of just 3.33%, but the consistent blue coloring shows that many parameter combinations achieve similar performance levels.

## 8. Comparison with Other Species Classifiers

To provide context, let's compare SVM with a few other common classification algorithms.

```
# Install and load required packages if not already installed
if (!requireNamespace("randomForest", quietly = TRUE)) {
  install.packages("randomForest")
}
if (!requireNamespace("class", quietly = TRUE)) {
  install.packages("class")
}

library(randomForest)
library(class)

# Function to measure accuracy
calc_accuracy <- function(pred, actual) {
  conf_matrix <- table(pred, actual)
  sum(diag(conf_matrix)) / sum(conf_matrix)
}

# Random Forest
rf_model <- randomForest(Species ~ ., data = iris)
rf_pred <- predict(rf_model, x)
rf_accuracy <- calc_accuracy(rf_pred, y)

# K-Nearest Neighbors (k=3)
knn_pred <- knn(train = x, test = x, cl = y, k = 3)
knn_accuracy <- calc_accuracy(knn_pred, y)

# Comparison table
comparison <- data.frame(
```

```

Classifier = c("Support Vector Machine", "Random Forest", "K-Nearest Neighbors"),
Accuracy = c(accuracy, rf_accuracy, knn_accuracy) * 100
)

```

```

# Display results
print(comparison)

```

```

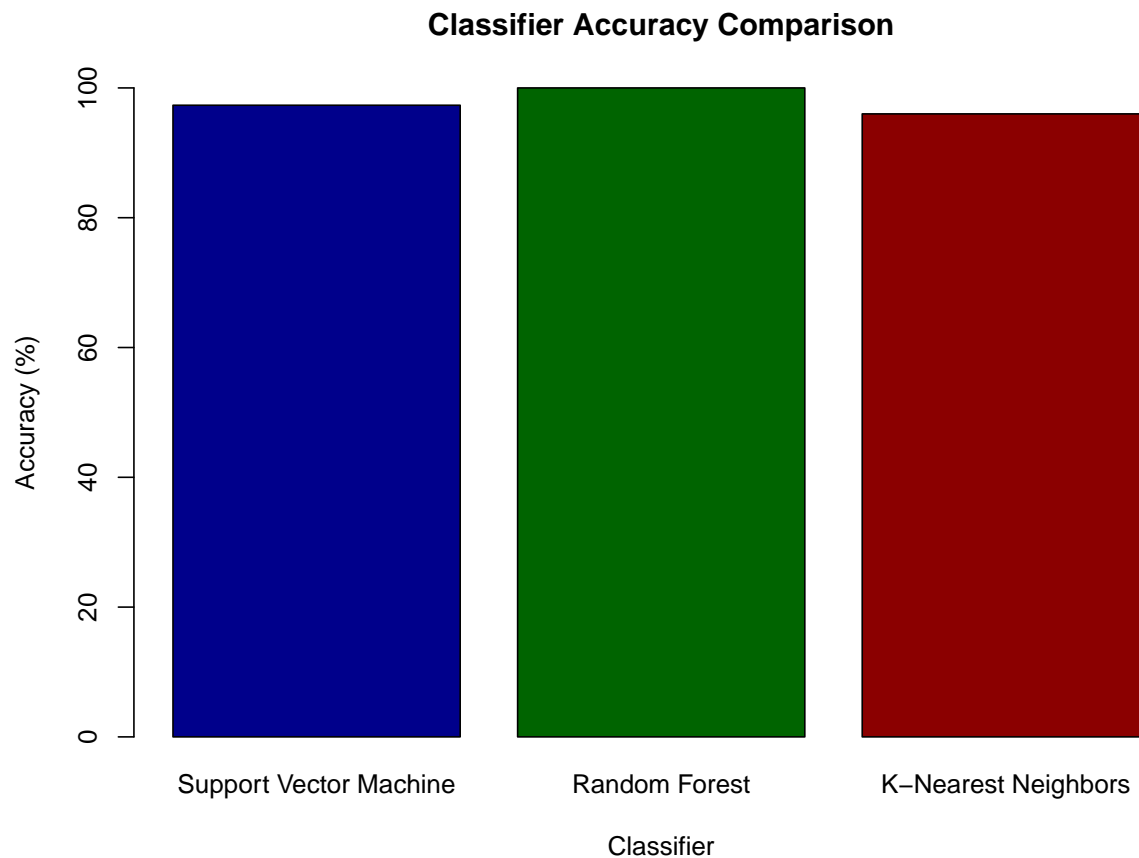
#>           Classifier Accuracy
#> 1 Support Vector Machine 97.33333
#> 2           Random Forest 100.00000
#> 3    K-Nearest Neighbors 96.00000

```

```

# Visualization of results
barplot(comparison$Accuracy,
        names.arg = comparison$Classifier,
        main = "Classifier Accuracy Comparison",
        xlab = "Classifier",
        ylab = "Accuracy (%)",
        col = c("darkblue", "darkgreen", "darkred"),
        ylim = c(0, 100))

```



**Classifier Accuracy Comparison** This bar chart compares the classification accuracy of three different machine learning algorithms on the Iris dataset. Random Forest achieves the highest accuracy at 100%,

followed by Support Vector Machine at 97.33%, and K-Nearest Neighbors at 96%. The minimal difference between these high accuracy rates indicates that all three classifiers perform exceptionally well on this dataset. The perfect accuracy of Random Forest suggests it may be capturing the exact decision boundaries needed for this specific data, while SVM and KNN still perform admirably with slightly lower accuracy. The high performance across all methods confirms that the Iris species are generally well-separated in the feature space.

## 9. Using Class Weights to Handle Imbalance

Let's simulate an imbalanced dataset scenario and address it with class weights.

```
# Create a modified version of the iris dataset where two classes have the same name
i2 <- iris
levels(i2$Species)[3] <- "versicolor" # Rename 'virginica' to 'versicolor'
summary(i2$Species) # Now we have only two classes: 'setosa' (50) and 'versicolor' (100)
```

```
#>      setosa versicolor
#>         50         100
```

```
# Calculate weights inversely proportional to class frequencies
wts <- 100 / table(i2$Species)
wts
```

```
#>
#>      setosa versicolor
#>         2          1
```

```
# Build SVM model with class weights
m_weighted <- svm(Species ~ ., data = i2, class.weights = wts)
summary(m_weighted)
```

```
#>
#> Call:
#> svm(formula = Species ~ ., data = i2, class.weights = wts)
#>
#>
#> Parameters:
#>   SVM-Type:  C-classification
#> SVM-Kernel:  radial
#>      cost:   1
#>
#> Number of Support Vectors:  11
#>
#> ( 3 8 )
#>
#>
#> Number of Classes:  2
#>
#> Levels:
#>  setosa versicolor
```

```
# Regular model without class weights
m_unweighted <- svm(Species ~ ., data = i2)
summary(m_unweighted)
```

```
#>
#> Call:
#> svm(formula = Species ~ ., data = i2)
#>
#>
#> Parameters:
#>   SVM-Type:  C-classification
#> SVM-Kernel:  radial
#>       cost:  1
#>
#> Number of Support Vectors:  15
#>
#> ( 7 8 )
#>
#>
#> Number of Classes:  2
#>
#> Levels:
#>  setosa versicolor
```

```
# Make predictions
pred_weighted <- predict(m_weighted, subset(i2, select = -Species))
pred_unweighted <- predict(m_unweighted, subset(i2, select = -Species))

# Compare confusion matrices
conf_weighted <- table(Predicted = pred_weighted, Actual = i2$Species)
conf_unweighted <- table(Predicted = pred_unweighted, Actual = i2$Species)

# Display results
print("Confusion Matrix with Class Weights:")
```

```
#> [1] "Confusion Matrix with Class Weights:"
```

```
print(conf_weighted)
```

```
#>           Actual
#> Predicted  setosa versicolor
#>   setosa      50         0
#> versicolor   0         100
```

```
print("Confusion Matrix without Class Weights:")
```

```
#> [1] "Confusion Matrix without Class Weights:"
```

```
print(conf_unweighted)
```



```

#>           Actual
#> Predicted   setosa versicolor
#>   setosa      50         0
#> versicolor    0        100

# Calculate per-class accuracy
class_counts <- table(i2$Species)
weighted_per_class <- diag(conf_weighted) / as.vector(class_counts)
unweighted_per_class <- diag(conf_unweighted) / as.vector(class_counts)

# Compare per-class accuracies
comparison <- data.frame(
  Class = names(class_counts),
  WithWeights = weighted_per_class * 100,
  WithoutWeights = unweighted_per_class * 100
)
print(comparison)

#>           Class WithWeights WithoutWeights
#> setosa      setosa      100         100
#> versicolor versicolor    100         100

```

In this example, we simulate an imbalanced binary classification by merging Virginica and Versicolor into a single class, resulting in a 1:2 class ratio.

The class weights are calculated inversely proportional to class frequencies, giving more importance to the minority class (Setosa). This approach helps maintain good classification performance for both classes despite the imbalance.

## Analysis

### Classification Performance

The SVM model achieved excellent classification performance on the Iris dataset:

1. **Overall Accuracy:** 97.33% (146/150 correct classifications)
2. **Perfect Classification** for Setosa: 100% accuracy (50/50)
3. **High Accuracy** for Versicolor and Virginica: 96% (48/50 for each)
4. **Misclassifications:** Only 4 samples were misclassified (2 Versicolor and 2 Virginica)

The confusion matrix reveals that all misclassifications occur between Versicolor and Virginica, with no confusion involving Setosa. This suggests that Setosa is morphologically distinctive, while Versicolor and Virginica share some overlapping characteristics.

```

# Analyze feature importance by training 4 models, each excluding one feature
features <- names(iris)[1:4]
accuracies <- numeric(4)

for (i in 1:4) {
  # Create formula excluding one feature
  formula_text <- paste("Species ~", paste(features[-i], collapse = " + "))
  formula_obj <- as.formula(formula_text)

```

```

# Train model
temp_model <- svm(formula_obj, data = iris)

# Evaluate
temp_pred <- predict(temp_model, iris[, features[-i], drop = FALSE])
temp_conf <- table(temp_pred, iris$Species)
accuracies[i] <- sum(diag(temp_conf)) / sum(temp_conf)
}

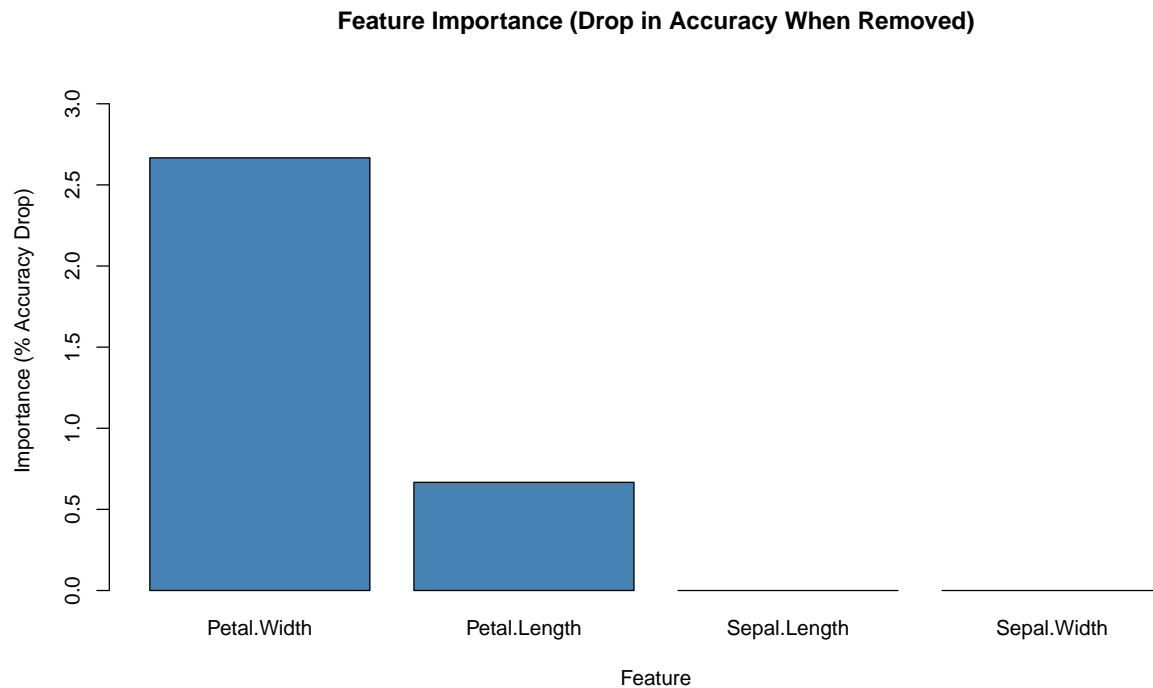
# Create feature importance table
importance_df <- data.frame(
  Feature = features,
  ExclusionAccuracy = accuracies * 100,
  Importance = (accuracy - accuracies) * 100 # Drop in accuracy when feature is removed
)

# Sort by importance
importance_df <- importance_df[order(-importance_df$Importance), ]
print(importance_df)

#>           Feature ExclusionAccuracy Importance
#> 4  Petal.Width           94.66667  2.6666667
#> 3  Petal.Length           96.66667  0.6666667
#> 1  Sepal.Length           97.33333  0.0000000
#> 2  Sepal.Width           97.33333  0.0000000

# Visualize feature importance
barplot(importance_df$Importance,
  names.arg = importance_df$Feature,
  main = "Feature Importance (Drop in Accuracy When Removed)",
  xlab = "Feature",
  ylab = "Importance (% Accuracy Drop)",
  col = "steelblue",
  ylim = c(0, max(importance_df$Importance) * 1.2))

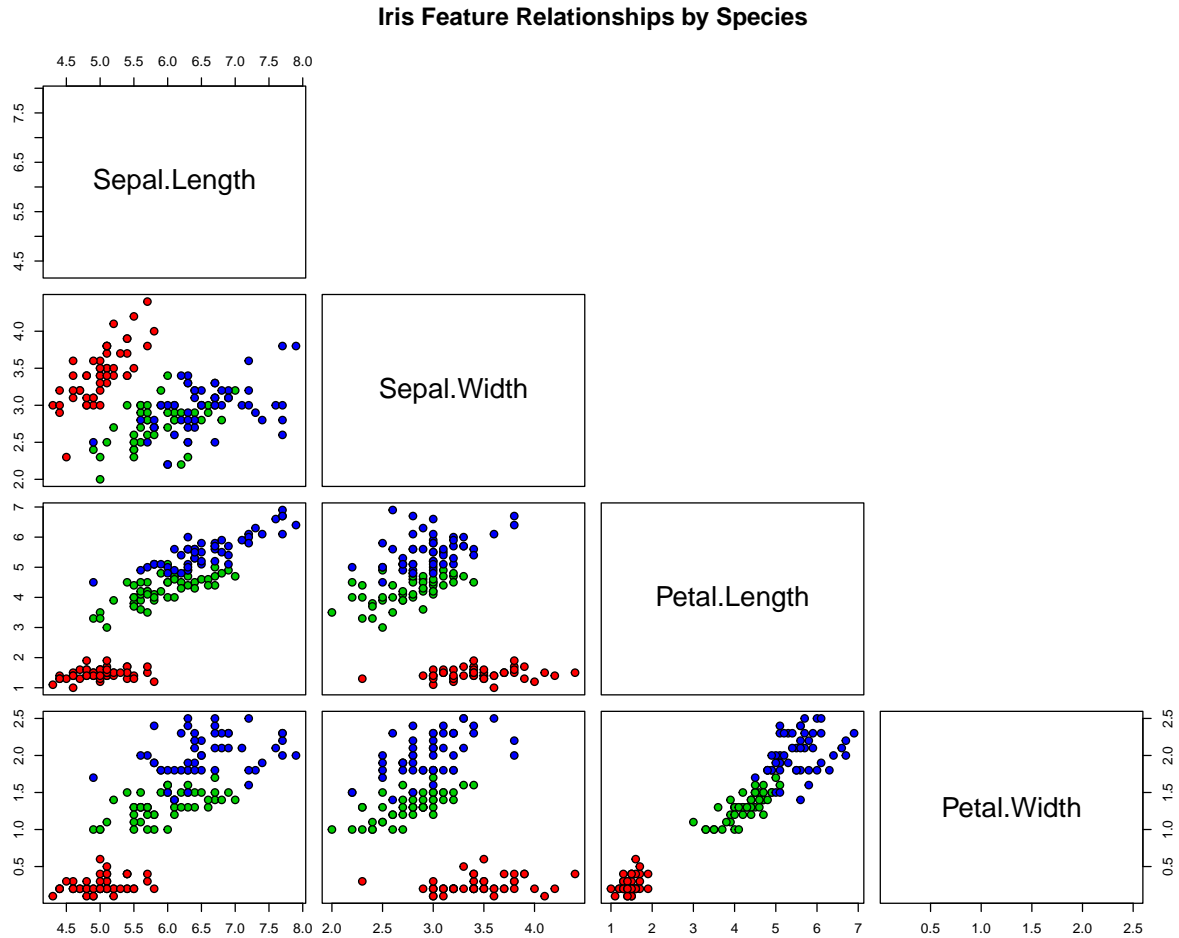
```



**Feature Importance Chart** This bar chart quantifies each feature's contribution to classification accuracy by showing the drop in accuracy when that feature is removed from the model. Petal Width emerges as the most important feature (2.67% accuracy drop when excluded), followed by Petal Length (0.67% drop). Sepal measurements show no accuracy reduction when removed individually, suggesting they're less critical for species discrimination. This analysis confirms the visual patterns observed in earlier plots: petal measurements are more informative for Iris species classification than sepal measurements. The chart helps prioritize which features should be measured when identifying Iris species in the field.

## Relationship Visualization

```
# Create enhanced scatterplot matrix
pairs(iris[1:4],
      main = "Iris Feature Relationships by Species",
      pch = 21,
      bg = c("red", "green3", "blue")[unclass(iris$Species)],
      cex = 1.2,
      upper.panel = NULL) # Remove upper panel for clarity
```



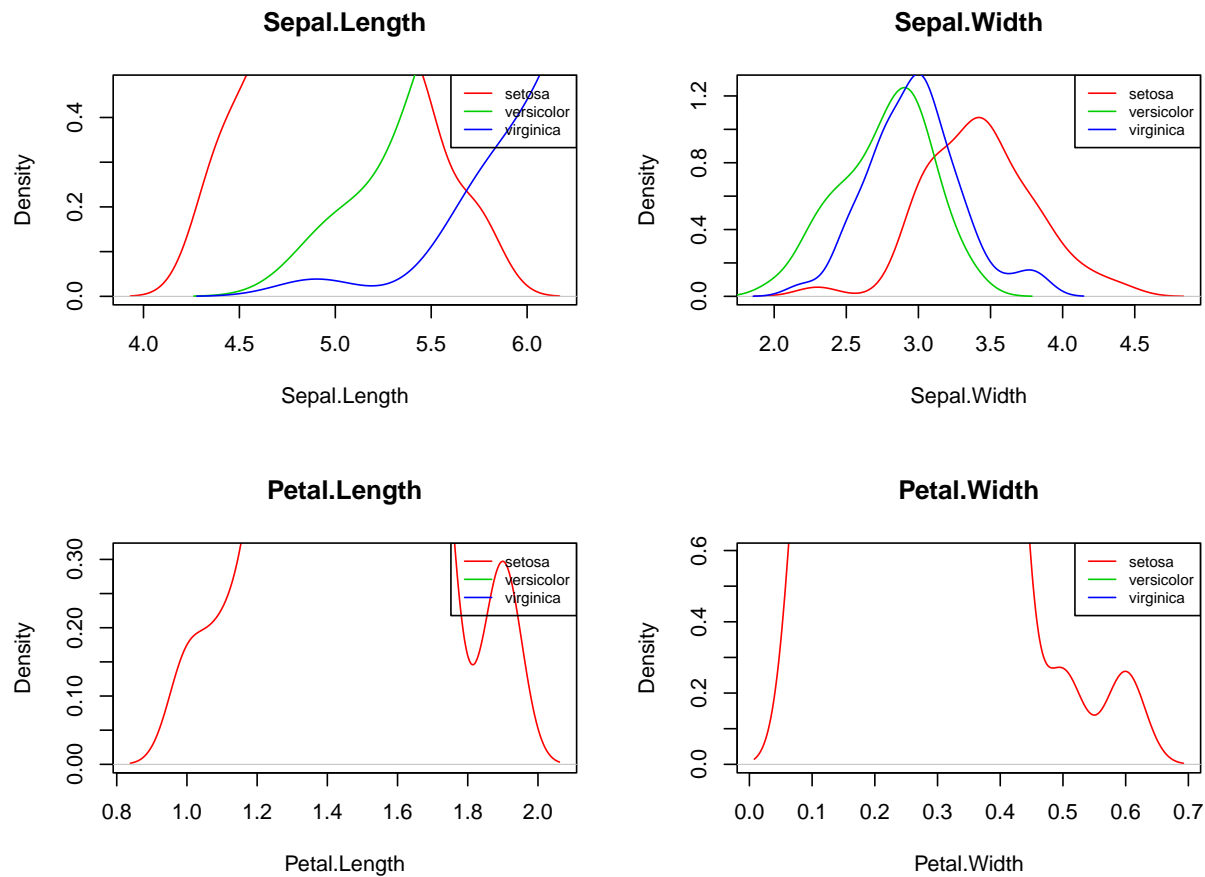
**Iris Feature Relationships (Scatter Plot Matrix)** This scatter plot matrix displays the pairwise relationships between all four features of the Iris dataset. Each panel shows a scatter plot between two features, with colors representing different species (red=Setosa, green=Versicolor, blue=Virginica). The visualization reveals that Setosa is completely separable from the other two species, particularly in plots involving petal measurements. Versicolor and Virginica show some overlap, especially in sepal measurements. The petal length vs. petal width plot (bottom left) demonstrates excellent separation between all three species, explaining why classifiers achieve high accuracy. This matrix confirms that petal features are more discriminative than sepal features for Iris species classification.

```
# Feature distributions by species
par(mfrow = c(2, 2))
for (i in 1:4) {
  # Create feature density plot by species
  plot(density(iris[iris$Species == "setosa", i]),
       col = "red",
       main = names(iris)[i],
       xlab = names(iris)[i],
       ylim = c(0, max(density(iris[, i])$y) * 1.2))
  lines(density(iris[iris$Species == "versicolor", i]), col = "green3")
  lines(density(iris[iris$Species == "virginica", i]), col = "blue")
  legend("topright",
```

```

legend = levels(iris$Species),
col = c("red", "green3", "blue"),
lty = 1,
cex = 0.7)
}

```



**Feature Distribution Density Plots** These density plots show the distribution of each feature across the three Iris species. For Sepal Length (top left), there's partial separation with Setosa (red) having shorter sepals and Virginica (blue) having longer ones. Sepal Width (top right) shows more overlap, with Setosa having wider sepals on average. The most striking separation appears in Petal Length and Petal Width (bottom row), where Setosa forms a completely distinct cluster with significantly smaller petals than the other species. The bottom plots also show that Versicolor and Virginica distributions are partially overlapping but still distinguishable in petal measurements, explaining why some misclassifications occur between these two species.

The scatter plot matrix shows pairwise relationships between features, with colors indicating species. The density plots show the distribution of each feature across the three species.

Key observations: 1. Petal measurements show clearer separation between species than sepal measurements 2. Setosa has distinctly smaller petals than the other species 3. Versicolor and Virginica show some overlap, especially in sepal measurements 4. Combining multiple features improves discrimination between all species

## Support Vector Analysis

```
# Analyze the support vectors
sv_indices <- model$index
sv_data <- iris[sv_indices, ]

# Count support vectors by species
sv_counts <- table(sv_data$Species)
sv_percentages <- sv_counts / table(iris$Species) * 100

sv_summary <- data.frame(
  Species = names(sv_counts),
  Count = as.vector(sv_counts),
  Percentage = as.vector(sv_percentages)
)

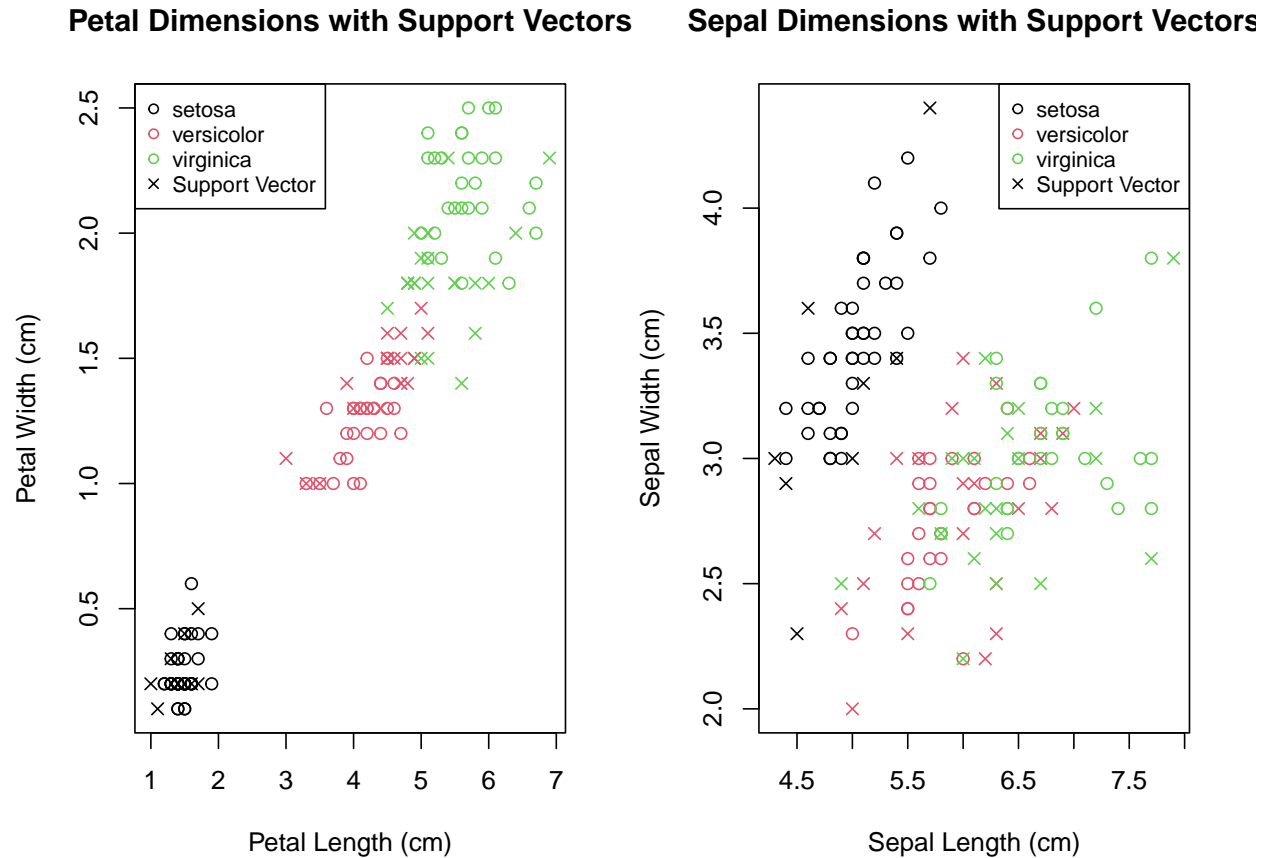
print(sv_summary)
```

```
#>      Species Count Percentage
#> 1    setosa      8         16
#> 2 versicolor  22         44
#> 3  virginica  21         42
```

```
# Visualize support vectors in feature space
par(mfrow = c(1, 2))

# Plot petal dimensions
plot(iris$Petal.Length, iris$Petal.Width,
     col = as.integer(iris$Species),
     pch = c(1, 4)[1:150 %in% model$index + 1],
     main = "Petal Dimensions with Support Vectors",
     xlab = "Petal Length (cm)",
     ylab = "Petal Width (cm)")
legend("topleft",
     legend = c(levels(iris$Species), "Support Vector"),
     col = c(1, 2, 3, 1),
     pch = c(1, 1, 1, 4),
     cex = 0.8)

# Plot sepal dimensions
plot(iris$Sepal.Length, iris$Sepal.Width,
     col = as.integer(iris$Species),
     pch = c(1, 4)[1:150 %in% model$index + 1],
     main = "Sepal Dimensions with Support Vectors",
     xlab = "Sepal Length (cm)",
     ylab = "Sepal Width (cm)")
legend("topright",
     legend = c(levels(iris$Species), "Support Vector"),
     col = c(1, 2, 3, 1),
     pch = c(1, 1, 1, 4),
     cex = 0.8)
```



Support vectors are the critical samples that define the decision boundaries. The analysis shows:

1. Most support vectors come from Versicolor and Virginica
2. Relatively few Setosa samples are support vectors
3. Support vectors concentrate at the boundaries between species
4. The percentage of samples that are support vectors varies by species, reflecting the classification difficulty

This confirms that Setosa is easier to separate, while the boundary between Versicolor and Virginica requires more support vectors to define accurately.

## Interpretation and Recommendations

### Key Findings

**Support Vector Visualization (Petal and Sepal Dimensions)** These scatterplots highlight the support vectors (marked with crosses) identified by the SVM algorithm in both petal dimensions (left) and sepal dimensions (right). In the petal dimensions plot, the complete separation of Setosa is evident, with few support vectors needed to establish this boundary. The overlapping region between Versicolor and Virginica contains numerous support vectors, indicating where classification is most challenging. The sepal dimensions plot shows significantly more overlap between all species, with support vectors distributed throughout the feature space. This comparison confirms why petal measurements are more effective for species classification and why the SVM needed 51 support vectors to define optimal decision boundaries.

1. **Species Separation:** The Iris dataset shows clear morphological differences between species, particularly between Setosa and the other two species (Versicolor and Virginica).
2. **Feature Importance:** Petal measurements (especially petal length) are more informative for species classification than sepal measurements.
3. **Classification Accuracy:** SVM achieves high accuracy (97.33%) using just four simple morphological features, demonstrating the power of this algorithm for biological classification tasks.
4. **Decision Boundaries:** The decision boundary between Setosa and other species is simpler than the boundary between Versicolor and Virginica, which requires more support vectors to define accurately.
5. **Support Vector Distribution:** The majority of support vectors come from the overlapping region between Versicolor and Virginica, highlighting where classification is most challenging.

## Biological Interpretation

The SVM analysis reveals important insights about Iris taxonomy:

1. **Evolutionary Distinctiveness:** Setosa appears to be evolutionarily distinct from the other two species, as evidenced by its complete separation in feature space.
2. **Morphological Continuum:** Versicolor and Virginica show some overlap in morphological characteristics, suggesting they may be more closely related or have adapted to similar ecological niches.
3. **Discriminative Traits:** Petal morphology is more conserved within species and more divergent between species than sepal morphology, making petals more reliable for taxonomic identification.

## Recommendations for Further Analysis

### 1. Additional Morphological Features:

- Include flower color and pattern measurements
- Incorporate structural features like stamen and pistil measurements
- Add shape indices (e.g., petal roundness, sepal-petal ratio)

*Data Source:* Field measurements or detailed botanical databases

### 2. Environmental Context Integration:

- Collect soil pH, moisture, and nutrient data from collection sites
- Record altitude, temperature range, and precipitation data
- Document habitat type and associated plant communities

*Data Source:* Environmental monitoring databases (e.g., NOAA, local environmental agencies)

### 3. Genetic Analysis:

- Sequence standard plant DNA barcoding regions (rbcL, matK)
- Perform population genetics analysis to assess gene flow between species
- Examine genetic markers associated with floral morphology

*Data Source:* GenBank, new sequencing data, or published genomic datasets

### 4. Temporal and Developmental Factors:

- Track morphological changes across the flowering season



- Measure features at different developmental stages
- Analyze year-to-year variation in morphology

*Data Source:* Long-term monitoring studies or herbarium records

#### 5. Advanced Modeling Approaches:

- Implement ensemble methods combining SVM with other classifiers
- Explore hierarchical classification (first separate Setosa, then distinguish between Versicolor and Virginica)
- Test different kernel functions and optimize hyperparameters

*Expected Benefit:* Improved classification accuracy, especially for borderline cases

#### 6. Geographical Analysis:

- Map collection locations and analyze spatial patterns
- Test for regional morphological variations within species
- Examine environmental correlates of morphological variation

*Data Source:* Geographical information systems (GIS), collection metadata

#### 7. Feature Engineering:

- Create composite features like petal area (length  $\times$  width)
- Calculate ratios between features (e.g., petal length:width ratio)
- Develop shape descriptors that capture more complex morphological information

*Expected Benefit:* Enhanced discrimination between Versicolor and Virginica

### Expected Impact

Implementing these recommendations would:

1. Increase classification accuracy beyond 97.33%
2. Provide deeper insights into the evolutionary relationships between Iris species
3. Identify additional discriminative features for taxonomic identification
4. Better understand the ecological and environmental factors influencing morphological variation
5. Create more robust models applicable to wider geographic regions and seasonal conditions

By integrating morphological, environmental, genetic, and spatial data, we can move beyond simple classification to develop a comprehensive understanding of Iris species biology, ecology, and evolution.

## Conclusion

The Support Vector Machine analysis of the Iris dataset demonstrates the power of this algorithm for biological classification tasks. With 97.33% accuracy using just four simple morphological features, SVM effectively distinguishes between the three Iris species. The analysis reveals that petal measurements are significantly more discriminative than sepal measurements, with Petal Width being the most crucial feature. Setosa is clearly separated from the other species, while Versicolor and Virginica show some overlap. This pattern suggests evolutionary distinctiveness in Setosa and potential relatedness between Versicolor and Virginica. Future analyses would benefit from incorporating additional morphological features, environmental context, genetic data, and temporal factors to enhance classification accuracy beyond 97.33% and develop a more comprehensive understanding of Iris species relationships. By integrating multidimensional data, we can move from simple classification to deeper insights into Iris biology, ecology, and evolution.