

# svmodt: An R Package for Linear SVM-Based Oblique Decision Trees

by Aneesh Agarwal, Jack Jewson, and Erik Sverdrup

**Abstract** An abstract of less than 150 words.

## 1 Introduction

## 2 Background

### Decision Trees

Decision Trees (DTs) are interpretable classification models that represent their decision-making process through a hierarchical, tree-like structure. This structure comprises internal nodes containing splitting criteria and terminal (leaf) nodes corresponding to class labels. The nodes are connected by directed edges, each representing a possible outcome of a splitting criterion. Formally, a DT can be expressed as a rooted, directed tree  $T = (G(V, E), v_1)$ , where  $V$  denotes the set of nodes,  $E$  represents the set of edges linking these nodes, and  $v_1$  is the root node.

If the tree  $T$  has  $m$  nodes, then for any  $j \in \{1, \dots, m\}$ , the set of child nodes of  $v_j \in V$  can be defined as:

$$N^+(v_j) = \{v_k \in V \mid k \in \{1, \dots, m\}, k \neq j, (v_j, v_k) \in E\}.$$

Here,  $N^+(v_j)$  denotes the set of nodes that are directly connected to  $v_j$  through outgoing edges, representing all possible child nodes that can be reached from  $v_j$  within the tree structure (Rivera-Lopez and Canul-Reich, 2018).

Decision tree algorithms can be categorized based on whether the same type of test is applied at all internal nodes. **Homogeneous trees** employ a single algorithm throughout (e.g., univariate or multivariate splits), whereas **hybrid trees** allow different algorithms such as linear discriminant functions,  $k$ -nearest neighbors, or univariate splits that can be used in different subtrees (Brodley and Utgoff, 1995). Hybrid trees exploit the principle of *selective superiority*, allowing subsets of the data to be modeled by the most appropriate classifier, thereby improving flexibility and accuracy.

### Univariate Decision Trees

Univariate Decision Trees (UDTs) trees represent axis-parallel hyperplanes dividing the instance space into several disjoint regions. Axis-parallel decision trees, such as CART and C4.5, represent two of the most widely used algorithms for classification tasks. The **CART (Classification and Regression Trees)** algorithm employs a binary recursive partitioning procedure capable of handling both continuous and categorical variables as predictors or targets. It operates directly on raw data without requiring binning. The tree is expanded recursively until no further splits are possible, after which **cost-complexity pruning** is applied to remove branches that contribute least to predictive performance. This pruning process generates a sequence of nested subtrees, from which the optimal model is selected using independent test data or cross-validation, rather than internal training measures (Breiman et al., 1984).

In contrast, **C4.5**, an extension of the earlier **ID3** algorithm (Quinlan, 1986), utilizing information theory measures such as **information gain** and **gain ratio** to select the most informative attribute for each split (Quinlan, 2014). C4.5 also includes mechanisms to handle missing attribute values by weighting instances according to the proportion of known data and employs an **error-based pruning** method to reduce overfitting. Although these techniques are effective across diverse datasets, studies have shown that the choice of pruning strategy and stopping criteria can significantly affect model performance across different domains (Mingers, 1989; Schaffer, 1992).

While UDTs are highly interpretability, they are characterised by several representational limitations. Such trees often grow unnecessarily large, as they must approximate complex relationship between features through multiple axis-aligned partitions. This can result in the replication of subtrees and repeated testing of the same feature along different paths, both of which reduce efficiency and hinder generalization performance (Pagallo and Haussler, 1990).

## Multivariate Decision Trees

Multivariate decision trees (MDTs) extends UDTs by allowing each internal node to perform splits based on linear or nonlinear combinations of multiple features. This flexibility enables the tree to form oblique decision boundaries that more accurately partition the instance space. For example, a single multivariate test such as  $x + y < 8$  can replace multiple univariate splits needed to approximate the same boundary. The construction of MDTs introduces several design considerations, including how to represent multivariate tests, determine their coefficients, select features to include, handle symbolic and missing data, and prune to avoid overfitting (Brodley and Utgoff, 1995).

Various optimization algorithms—such as recursive least squares (Young, 1984), the pocket algorithm (Gallant, 1986), or thermal training (Frean, 1990)—may be used to estimate the weights. However, MDTs trade interpretability for representational power and often require additional mechanisms for **local feature selection**, such as *sequential forward selection* (SFS) or *sequential backward elimination* (SBE) (Kittler, 1986).

Empirical comparisons across multiple datasets demonstrate that multivariate trees generally achieve higher accuracy and smaller tree sizes than their univariate counterparts, though this comes at the cost of reduced interpretability. Moreover, MDTs retain key advantages of standard decision trees—such as sequential split criteria evaluation and transparent decision procedures—while offering improved modeling flexibility for complex datasets (Kozioł and Wozniak, 2009; Friedl and Brodley, 1997; Liu and Setiono, 1998; Cañete et al., 2021).

## Support Vector Machines (SVMs)

Support Vector Machines (SVMs) are powerful supervised learning models used for classification and regression tasks. They aim to determine an optimal separating hyperplane that maximizes the margin between different classes in the data. This margin-based approach enhances the generalization ability of the model, making SVMs robust and effective for many real-world problems (Cristianini, 2000).

### Linear SVMs

A simplest **linear SVMs** construct a separating hyperplane in an  $n$ -dimensional space such that the margin between the classes is maximized. Given a training dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , where  $\mathbf{x}_i \in \mathbb{R}^n$  and  $y_i \in \{-1, +1\}$ , the decision function is defined as:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b),$$

where  $\mathbf{w}$  is the weight vector perpendicular to the hyperplane, and  $b$  is the bias term. The optimal hyperplane is the one that maximizes the distance between the closest points of each class (the **support vectors**) and the hyperplane itself (Cortes and Vapnik, 1995).

Mathematically, the optimization problem for a hard-margin SVM (i.e., assuming the data are linearly separable) can be formulated as:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N.$$

Here,  $\|\mathbf{w}\|$  represents the norm of the weight vector and acts as a regularization term that controls the complexity of the model. The constraint ensures that all data points are correctly classified and lie outside the margin boundaries.

However, in most practical situations, perfect linear separability is not achievable. To address this, **soft-margin SVMs** introduce slack variables  $\xi_i \geq 0$  to allow certain violations of the margin constraints, resulting in the following optimization problem:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

subject to:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, 2, \dots, N.$$

The parameter  $C > 0$  controls the trade-off between maximizing the margin and minimizing the classification error on the training data. A large  $C$  penalizes mis-classifications heavily, leading to a

**Table 1:** Commonly used kernel functions and their parameters.

| Kernel Type | Mathematical Definition   | Key Parameters     |
|-------------|---|--------------------|
| Linear      | $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$  | None               |
| Polynomial  | $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + 1)^d$                                | Degree $d$         |
| Gaussian    | $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right)$ | Bandwidth $\sigma$ |
| RBF         | $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma\ \mathbf{x}_i - \mathbf{x}_j\ ^2)$                        | $\gamma$           |
| Sigmoid     | $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\kappa \mathbf{x}_i^\top \mathbf{x}_j + \theta)$                 | $\kappa, \theta$   |

narrower margin, whereas a smaller  $C$  allows more flexibility, potentially improving generalization in the presence of noise.

The solution to this constrained optimization problem is obtained using **Lagrange multipliers**, resulting in a dual formulation expressed as:

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

subject to:

$$\sum_{i=1}^N \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C.$$

The data points corresponding to non-zero  $\alpha_i$  values are the **support vectors**, which define the decision boundary. The resulting decision function for a new observation  $\mathbf{x}$  is given by:

$$f(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^N \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b \right).$$

This formulation highlights one of the most important properties of SVMs — the decision boundary depends only on a subset of the training data (the support vectors), making SVMs both efficient and robust in representing the learned model (Cervantes et al., 2020).

## Non-Linear SVMs

While linear classifiers provide useful insights, they are often inadequate for real-world datasets, where classes are not linearly separable. In such cases, SVMs can be extended to create **nonlinear decision boundaries** by mapping the input vectors into a higher-dimensional **feature space** using a nonlinear transformation  $\phi: \mathbb{R}^n \rightarrow \mathcal{F}$ . The linear transformation is then achieved in the transformed space using:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \phi(\mathbf{x}) + b).$$

However, directly computing  $\phi(\mathbf{x})$  can be computationally expensive. To address this, SVMs employ the **kernel trick**, where the dot product in the feature space is replaced with a kernel function  $K(\mathbf{x}_i, \mathbf{x}_j)$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle.$$

The resulting decision function becomes:

$$f(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right),$$

where  $\alpha_i$  are the Lagrange multipliers obtained during training.

For a function  $K$  to be a valid kernel, it must satisfy **Mercer's condition** (Vapnik, 2013) i.e., the kernel matrix must be symmetric and positive semi-definite. Table 1

Although SVMs exhibit strong theoretical foundations and robust generalization capabilities, they present several practical limitations. Model performance is highly dependent on the appropriate selection of hyperparameters such as the regularization term ( $C$ ) and kernel parameters (e.g.,  $\gamma$ ), which govern the trade-off between margin maximization and misclassification tolerance (Nanda et al., 2018). Training an SVM requires solving a quadratic programming (QP) optimization problem involving an  $n \times n$  kernel matrix, where  $n$  denotes the number of training samples, leading to quadratic growth in both computational time and memory usage (Dong et al., 2005). This makes SVMs computationally

expensive for large-scale datasets. Moreover, SVMs are inherently designed for binary classification, necessitating decomposition strategies such as One-vs-One and One-vs-All for multi-class problems (Hsu and Lin, 2002). Their performance also tends to degrade in imbalanced data settings, where the decision boundary becomes biased toward the majority class (Cervantes et al., 2020).

### Hybrid Decision Trees with Support Vector Machines

The earliest formalization of Support Vector Machine–based Decision Trees (DTSVMs) was introduced by Bennett and Blue (1998), who extended the principles of Statistical Learning Theory (Vapnik, 2013) and Structural Risk Minimization (SRM) to the construction of binary decision trees. Each node of a decision tree was treated as an SVM that partitions data along an optimal hyperplane. This approach allowed for multivariate SVM decisions at each node. TAKAHASHI and Abe (2002) extended DTSVMs to multiclass problems, using a recursive tree-based partitioning where the root separates the most separable class (or classes) from the rest. Subsequent nodes repeat this process until each leaf contains a single class, ensuring that all regions of the feature space are classified.

Subsequent studies have built upon this foundation to address scalability, optimization, and generalization issues. Optimal Decision Tree SVM (ODT-SVM) (Bala and Agrawal, 2011) introduced split-selection criteria based on Gini index, information gain, and scatter matrix separability to balance tree interpretability and margin-based precision.

### Oblique and Rotation-Based Extensions

Recent research has explored oblique and rotation-based decision ensembles that generalize the DTSVM concept. Oblique Double Random Forests with MPSVM (MPDRaF) (Ganaie et al., 2022) introduce multivariate (oblique) splits at each node using Multi-Plane SVM (MPSVM) formulations (Mangasarian and Wild, 2006), enhancing the geometric flexibility of decision boundaries. Regularization variants—such as Tikhonov, axis-parallel, and null-space—address sample-size limitations at deeper nodes, ensuring better generalization.

Similarly, Rotation-based Double Random Forests (DRaF) (Ganaie et al., 2022) employ Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) transformations at non-terminal nodes to create diverse subspaces, improving ensemble diversity and classification stability.

### Modern Multi-class Integration

Recent works such as the Oblique Decision Tree Ensemble (ODTE) (Montañana et al., 2025) embed multiple SVM classifiers (e.g., one-vs-one or one-vs-rest) directly at each split, dynamically selecting the model that minimizes class impurity. This enables efficient n-ary classification within a single tree structure, mitigating the scalability and imbalance limitations of traditional binary SVM extensions.

### Related Software Packages

Several software packages provide functionality related to SVM-based decision trees, but fully-featured implementations of SVM Oblique Decision Trees (SVMODT) remain limited, particularly in R. In R, very scarce SVMODT implementation is available. The **PPforest** package implements projection pursuit random forests, where trees are constructed by splitting on linear combinations of variables to maximize class separation, supporting multi-class problems and variable importance measures. In Python, **STree** (Montañana et al., 2021) provides an oblique decision tree classifier based on SVM nodes, where each node is built and split using scikit-learn's **SVC** models (Pedregosa et al., 2011). Java-based frameworks such as **Weka** provide classical decision trees (**J48**) and SVMs (**SMO**) that can be adapted for hierarchical SVM trees (Menkovski et al., 2008). Overall, while multiple languages provide the building blocks—SVM training, decision tree construction, and ensemble methods—there is currently no widely available, fully integrated SVMODT software in R, highlighting the need for custom implementations or adaptation of existing tools for research and applied purposes.

## 3 Methodology

To better understand the structure and mechanics of decision trees, we first implemented a standard axis-parallel decision tree in R from scratch before extending it into the oblique (SVM-based) variant.

## Custom Univariate Decision Tree with Integrated Cost-Complexity Pruning

This implementation presents a decision tree algorithm with **integrated cost-complexity pruning**, where pruning decisions are made during tree construction rather than as a post-processing step. The algorithm supports multiple splitting criteria and performs real-time pruning based on cost-complexity analysis.

### Algorithm Description

At each node  $v$  with samples  $\mathcal{D}_v$  and labels  $\mathcal{Y}_v$ , the algorithm:

1. **Compute Node Statistics:** At each node of the decision tree, several key statistics are computed to guide the splitting process. The class probabilities are calculated as  $p_c = \frac{n_c}{n}$ , where  $n_c$  is the count of samples belonging to class  $c$  in the node, and  $n = |\mathcal{D}_v|$  is the total number of samples at the node. The predicted class is assigned as the class with the highest probability,

$$\hat{y} = \arg \max_c p_c.$$

2. **Select Optimal Split:** The optimal split at each node is determined using one of three criteria:

- **Gini Impurity** is defined as

$$\text{Gini}(\mathcal{Y}) = 1 - \sum_{c=1}^K p_c^2,$$

where  $p_c$  is the proportion of samples belonging to class  $c$  in node  $\mathcal{Y}$  (Breiman et al., 1984).

- **Information Gain** evaluates the reduction in node impurity and is computed as

$$\text{IG}(f, s) = H(\mathcal{Y}_v) - \sum_{i \in \{L, R\}} \frac{n_i}{n} H(\mathcal{Y}_i),$$

where  $H(\mathcal{Y})$  denotes the impurity of a node (entropy or Gini),  $n_i$  is the number of samples in the child node, and  $n$  is the number of samples in the parent node [Breiman et al. (1984)].

- **Gain Ratio** further adjusts Information Gain by penalizing splits that create many small partitions (Quinlan, 2014) and is defined as

$$\text{GR}(f, s) = \frac{\text{IG}(f, s)}{\text{SI}(f, s)},$$

where the split information is

$$\text{SI}(f, s) = - \sum_{i \in \{L, R\}} \frac{n_i}{n} \log_2 \left( \frac{n_i}{n} \right).$$

3. **Evaluates Pruning Criteria:** To ensure the tree remains interpret-able and avoids over-fitting, integrated pruning is performed using cost-complexity analysis. The cost of a sub-tree  $T_v$  is given by

$$R_\alpha(T_v) = R(T_v) + \alpha \|T_v\|,$$

where  $R(T_v)$  is the miss-classification error,  $|T_v|$  is the number of leaves in the sub-tree, and  $\alpha \geq 0$  is a complexity parameter controlling the trade-off between accuracy and simplicity. For a single leaf node  $t_v$ , the cost is

$$R_\alpha(t_v) = R(t_v) + \alpha.$$

A sub-tree  $T_v$  is pruned, i.e., replaced by a single leaf node, if

$$R_\alpha(t_v) \leq R_\alpha(T_v),$$

ensuring that further splits are retained only when they yield sufficient reduction in classification error to justify the added complexity.

The custom uni variate decision tree implementation in R combines axis-parallel splits with integrated cost-complexity pruning, offering several key differences from standard CART/[rpart](#) trees. Numeric features are split only at midpoints where class labels change, while categorical features use binary splits at each level, reducing computational complexity. Splitting criteria include Gini impurity, information gain, and Gain Ratio—the latter penalizing splits with high split information to favor balanced partitions. Integrated pruning is applied during tree construction using a fixed  $\alpha$  parameter, replacing sub-trees with leaves whenever the cost-complexity criterion is satisfied, unlike CART/[rpart](#),

**Algorithm 1** Decision Tree with Integrated Cost-Complexity Pruning

---

```

1: Procedure GenerateTree( $\mathcal{Y}, \mathbf{X}, d, d_{\max}, \alpha$ )
2:  $n \leftarrow |\mathcal{Y}|$ 
3:  $\mathbf{p} \leftarrow \text{ClassProbabilities}(\mathcal{Y})$ 
4:  $\hat{y} \leftarrow \arg \max_c p_c$ 
5: if  $d \geq d_{\max}$  or  $|\text{unique}(\mathcal{Y})| \leq 1$  then
6:   return LeafNode( $\hat{y}, \mathbf{p}, n$ )
7: end if
8:  $(f^*, s^*, \text{score}) \leftarrow \text{FeatureSelector}(\mathcal{Y}, \mathbf{X}, \text{criterion})$ 
9: if  $f^* = \text{null}$  or  $\text{score} = -\infty$  then
10:  return LeafNode( $\hat{y}, \mathbf{p}, n$ )
11: end if
12:  $\mathcal{I}_L, \mathcal{I}_R \leftarrow \text{Split}(\mathbf{X}[f^*], s^*)$ 
13: if  $|\mathcal{I}_L| = 0$  or  $|\mathcal{I}_R| = 0$  then
14:  return LeafNode( $\hat{y}, \mathbf{p}, n$ )
15: end if
16:  $T_L \leftarrow \text{GENERATE TREE}(\mathcal{Y}[\mathcal{I}_L], \mathbf{X}[\mathcal{I}_L], d + 1, d_{\max}, \alpha)$ 
17:  $T_R \leftarrow \text{GENERATE TREE}(\mathcal{Y}[\mathcal{I}_R], \mathbf{X}[\mathcal{I}_R], d + 1, d_{\max}, \alpha)$ 
18:  $T_v \leftarrow \text{InternalNode}(f^*, s^*, T_L, T_R, \hat{y}, \mathbf{p}, n)$ 
19: {Integrated pruning decision}
20:  $R_{\text{subtree}} \leftarrow \text{SubtreeError}(T_v) + \alpha \cdot \text{CountLeaves}(T_v)$ 
21:  $R_{\text{leaf}} \leftarrow (1 - \max(\mathbf{p})) + \alpha$ 
22: if  $R_{\text{leaf}} \leq R_{\text{subtree}}$  then
23:  return LeafNode( $\hat{y}, \mathbf{p}, n$ ) {Prune subtree}
24: else
25:  return  $T_v$ 
26: end if
27: End Procedure

```

---

which performs post-pruning via cross-validation. This single-pass pruning approach reduces training time and space requirements, while still producing interpret-able trees. Predictions follow standard tree traversal, returning majority-class labels or probability estimates at leaves.

## Support Vector Machine based Oblique Decision Trees

### Algorithm Description

The SVM Oblique Decision Tree (SVMODT) constructs a binary classification decision tree where each internal node  $v$  at depth  $d$ . The algorithm is implemented in R using the [e1071](#) performs the following operations:

1. **Feature Selection:** A subset of  $m_d$  features is dynamically selected from the available feature set  $\mathcal{F}$ . This selection can be done using one of several strategies, including random selection, sampling, mutual information ranking (using the [FSelectorRcpp](#) package), or correlation-based selection. This allows the tree to adaptively focus on the most informative features at each node.
2. **Feature Scaling:** Once the features are selected, z-score normalization is applied to standardize them. For the selected feature matrix  $\mathbf{X}_v \in \mathbb{R}^{n \times m_d}$ , the scaled features are computed as

$$\mathbf{X}_v^{\text{scaled}} = \frac{\mathbf{X}_v - \boldsymbol{\mu}_v}{\sigma_v},$$

where  $\boldsymbol{\mu}_v$  and  $\sigma_v$  denote the mean and standard deviation of the features in node  $v$ .

3. **SVM Training:** A linear SVM is then trained on the scaled features. Optional class weights  $w_c$  can be applied to handle imbalanced data. The SVM optimization problem is formulated as

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n w_{y_i} \xi_i,$$

where  $C$  is the regularization parameter,  $\xi_i$  are slack variables, and  $w_{y_i}$  represents the class-

specific weight for sample  $i$ .

4. **Node Splitting:** After training, the decision values for each sample are computed as  $f(x_i) = \mathbf{w}^T x_i + b$ . Samples are then partitioned into left and right child nodes according to the sign of  $f(x_i)$ :
  - Left child:  $\{i : f(x_i) > 0\}$
  - Right child:  $\{i : f(x_i) \leq 0\}$

This process recursively continues for each child node until stopping criteria are met.

---

**Algorithm 2** SVM Oblique Decision Tree Construction
 

---

```

1: Procedure SVMSPplit( $\mathcal{D}, d, d_{\max}, n_{\min}$ )
2:  $n \leftarrow |\mathcal{D}|, \mathcal{Y} \leftarrow \text{labels}(\mathcal{D})$ 
3: if  $d > d_{\max}$  or  $|\text{unique}(\mathcal{Y})| = 1$  or  $n < n_{\min}$  then
4:   return LeafNode( $\mathcal{Y}$ ) {Stopping criteria}
5: end if
6:  $m_d \leftarrow \text{DynamicMaxFeatures}(d, m_{\text{base}}, \text{strategy})$ 
7:  $\mathcal{F}_d \leftarrow \text{SelectFeatures}(\mathcal{D}, m_d, \text{method})$ 
8:  $\mathbf{X}_{\text{scaled}}, \text{scaler} \leftarrow \text{Scale}(\mathcal{D}[\mathcal{F}_d])$ 
9: if  $|\mathcal{F}_d| = 0$  then
10:  return LeafNode( $\mathcal{Y}$ ) {No valid features}
11: end if
12:  $w_c \leftarrow \text{CalculateClassWeights}(\mathcal{Y}, \text{strategy})$ 
13:  $\text{SVM} \leftarrow \text{FitLinearSVM}(\mathbf{X}_{\text{scaled}}, \mathcal{Y}, w_c)$ 
14: if  $\text{SVM} = \text{null}$  then
15:  return LeafNode( $\mathcal{Y}$ ) {SVM fitting failed}
16: end if
17:  $\mathbf{f} \leftarrow \text{SVM.DecisionValues}(\mathbf{X}_{\text{scaled}})$ 
18:  $\mathcal{I}_L \leftarrow \{i : f_i > 0\}, \mathcal{I}_R \leftarrow \{i : f_i \leq 0\}$ 
19: if  $|\mathcal{I}_L| = 0$  or  $|\mathcal{I}_R| = 0$  then
20:  return LeafNode( $\mathcal{Y}$ ) {Ineffective split}
21: end if
22:  $\text{left} \leftarrow \text{SVMSPplit}(\mathcal{D}[\mathcal{I}_L], d + 1, d_{\max}, n_{\min})$ 
23:  $\text{right} \leftarrow \text{SVMSPplit}(\mathcal{D}[\mathcal{I}_R], d + 1, d_{\max}, n_{\min})$ 
24: return InternalNode( $\text{SVM}, \mathcal{F}_d, \text{scaler}, \text{left}, \text{right}$ )
25: End Procedure
  
```

---

**Key Hyper-parameters**

1. **Feature Selection:** The algorithm supports three strategies for selecting features at each node. The maximum number of features considered at depth  $d$  is defined as:

$$m_d = \begin{cases} m_{\text{base}} & \text{constant strategy} \\ \lfloor m_{\text{base}} \cdot \alpha^{d-1} \rfloor & \text{decrease strategy} \\ \text{Uniform}(\lfloor p \cdot \ell_{\min} \rfloor, \lfloor p \cdot \ell_{\max} \rfloor) & \text{random strategy} \end{cases}$$

Here,  $\alpha \in (0, 1]$  is the decrease rate for the “decrease strategy,”  $p$  is the total number of features,  $\ell_{\min}$  and  $\ell_{\max}$  define the fractional bounds for the random strategy, and  $m_{\text{base}}$  is the base number of features at the first depth.

2. **Feature Penalty:** To encourage diversity in feature usage across the tree, previously used features have their selection weight reduced by a factor of  $(1 - \lambda)$ , where  $\lambda \in [0, 1)$ .
3. **Class Weight Handling:** To address class imbalance, the algorithm allows four weighting schemes:
  - **None:**  $w_c = 1$  for all classes.
  - **Balanced:**  $w_c = \frac{n}{K \cdot n_c}$ , where  $K$  is the total number of classes and  $n_c$  is the number of samples in class  $c$ .
  - **Balanced sub-sample:**  $w_c = \frac{1}{n_c} \cdot \frac{K}{\sum_{c'} 1/n_{c'}}$ , which adjusts weights for sub-sampled data.
  - **Custom:** Users can define their own class weights.



## Interaction with SVM Training and Node Splitting

The hyper-parameters described above directly influence the behavior of the SVM-based decision tree at each node. The feature selection strategy determines which subset of features  $\mathcal{F}_d$  is used to fit the linear SVM at depth  $d$ , which affects the orientation and effectiveness of the decision hyper-plane. Feature penalties ensure that no single feature dominates multiple splits, promoting diversity in the learned splits and improving generalization. Once the SVM is trained on the selected and scaled features, the decision values are used to partition the samples into left and right child nodes. By controlling the number of features, penalizing repeated use, and adjusting class weights, the tree can achieve a balance between predictive accuracy, interpret-ability, and computational efficiency.

## 4 Data

In this study, we evaluate the SVM-based oblique decision tree (SVMODT) on two widely used benchmark datasets: **Palmer Penguins** and **Wisconsin Diagnostic Breast Cancer (WDBC)**. These datasets provide a mix of multi-class and binary classification problems, as well as numeric and categorical features, making them suitable for testing the flexibility and interpretability of SVMODTs.

### Palmer Penguins

The `palmerpenguins` data set is a multi-class data set containing morphological measurements for three penguin species (*Adelie*, *Chinstrap*, and *Gentoo*) observed on the Palmer Archipelago, Antarctica. The data set includes 344 complete observations with the following features:

- **bill\_length\_mm**: Length of the bill (numeric)
- **bill\_depth\_mm**: Depth of the bill (numeric)
- **flipper\_length\_mm**: Length of the flipper (numeric)
- **body\_mass\_g**: Body mass (numeric)
- **sex**: Sex of the penguin (categorical)
- **species**: Target class (categorical with three levels)

### Wisconsin Diagnostic Breast Cancer (WDBC)

The **WDBC** (Street et al., 1993) data set is a binary classification data set derived from fine needle aspirates of breast tissue. It contains 569 observations with 30 numeric features computed from digitized images of cell nuclei. The features include measures such as **radius**, **texture**, **perimeter**, **area**, **smoothness**, **compactness**, **concavity**, **concave points**, **symmetry**, **fractal dimension**. The target variable is **diagnosis** (Malignant (M) or benign (B) tumor). This data set is widely used for testing high-dimensional binary classification algorithms. Its numeric nature allows direct evaluation of SVM splits and the effect of class weighting strategies on imbalanced classes.

### Data Preparation

For both data sets, the following pre-processing steps are applied before training the SVMODT:

1. **Missing Values**: Rows with missing values are removed.
2. **Class Filtering**: For Palmer Penguins, rows corresponding to the *Gentoo* species are removed using `dplyr` to create a binary classification problem.
3. **Categorical Encoding**: Categorical variables (e.g., **sex** in Palmer Penguins) are encoded as factor levels suitable for training.

These pre-processing steps ensure that SVMODT can learn meaningful oblique splits across both data sets while maintaining interpret-ability.



## 5 Code Examples

## 6 Results

## 7 Discussion

## Bibliography

- M. Bala and R. Agrawal. Optimal decision tree based multi-class support vector machine. *Informatica*, 35(2), 2011. [p4]
- K. Bennett and J. Blue. A support vector machine approach to decision trees. In *1998 IEEE International Joint Conference on Neural Networks Proceedings. IEEE World Congress on Computational Intelligence (Cat. No.98CH36227)*, volume 3, pages 2396–2401 vol.3, 1998. doi: 10.1109/IJCNN.1998.687237. [p4]
- L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth, 1984. URL <http://lyle.smu.edu/~mhd/8331f06/cart.pdf>. [p1, 5]
- C. E. Brodley and P. E. Utgoff. Multivariate decision trees. *Machine Learning*, 19:45–77, 1995. URL <https://api.semanticscholar.org/CorpusID:16836572>. [p1, 2]
- L. Cañete, R. Monroy, and M. Medina-Pérez. A review and experimental comparison of multivariate decision trees. *IEEE Access*, PP:1–1, 08 2021. doi: 10.1109/ACCESS.2021.3102239. [p2]
- J. Cervantes, F. Garcia-Lamont, L. Rodríguez-Mazahua, and A. Lopez. A comprehensive survey on support vector machine classification: Applications, challenges and trends. *Neurocomputing*, 408:189–215, 2020. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2019.10.118>. URL <https://www.sciencedirect.com/science/article/pii/S0925231220307153>. [p3, 4]
- C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995. [p2]
- N. Cristianini. An introduction to support vector machines and other kernel-based learning methods, 2000. [p2]
- J.-x. Dong, A. Krzyzak, and C. Y. Suen. Fast svm training algorithm with decomposition on very large data sets. *IEEE transactions on pattern analysis and machine intelligence*, 27(4):603–618, 2005. [p3]
- M. R. Frean. *Small nets and short paths: Optimising neural computation*. PhD thesis, 1990. [p2]
- M. Friedl and C. Brodley. Decision tree classification of land cover from remotely sensed data. *Remote Sensing of Environment*, 61(3):399–409, 1997. ISSN 0034-4257. doi: [https://doi.org/10.1016/S0034-4257\(97\)00049-7](https://doi.org/10.1016/S0034-4257(97)00049-7). URL <https://www.sciencedirect.com/science/article/pii/S0034425797000497>. [p2]
- S. I. Gallant. Optimal linear discriminants. *Eighth International Conference on Pattern Recognition*, pages 849–852, 1986. URL <https://cir.nii.ac.jp/crid/1573668924476144256>. [p2]
- M. Ganaie, M. Tanveer, P. Suganthan, and V. Snasel. Oblique and rotation double random forest. *Neural Networks*, 153:496–517, 2022. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2022.06.012>. URL <https://www.sciencedirect.com/science/article/pii/S0893608022002258>. [p4]
- C.-W. Hsu and C.-J. Lin. A comparison of methods for multiclass support vector machines. *IEEE transactions on Neural Networks*, 13(2):415–425, 2002. [p4]
- J. Kittler. Feature selection and extraction. *Handbook of Pattern Recognition and Image Processing*, pages 59–83, 1986. URL <https://cir.nii.ac.jp/crid/1573950400671608448>. [p2]
- M. Koziol and M. Wozniak. *Multivariate Decision Trees vs. Univariate Ones*, pages 275–284. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009. ISBN 978-3-540-93905-4. doi: 10.1007/978-3-540-93905-4\_33. URL [https://doi.org/10.1007/978-3-540-93905-4\\_33](https://doi.org/10.1007/978-3-540-93905-4_33). [p2]
- H. Liu and R. Setiono. Feature transformation and multivariate decision tree induction. In S. Arikawa and H. Motoda, editors, *Discovery Science*, pages 279–291, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg. ISBN 978-3-540-49292-4. [p2]

- O. Mangasarian and E. Wild. Multisurface proximal support vector machine classification via generalized eigenvalues. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(1):69–74, 2006. doi: 10.1109/TPAMI.2006.17. [p4]
- V. Menkovski, I. T. Christou, and S. Efremidis. Oblique decision trees using embedded support vector machines in classifier ensembles. In *2008 7th IEEE International Conference on Cybernetic Intelligent Systems*, pages 1–6, 2008. doi: 10.1109/UKRICIS.2008.4798937. [p4]
- J. Mingers. An empirical comparison of pruning methods for decision tree induction. *Machine Learning*, 4:227–243, 01 1989. doi: 10.1023/A:1022604100933. [p1]
- R. Montañana, J. A. Gámez, and J. M. Puerta. Stree: a single multi-class oblique decision tree based on support vector machines, 2021. [p4]
- R. Montañana, J. A. Gámez, and J. M. Puerta. Odte—an ensemble of multi-class svm-based oblique decision trees. *Expert Systems with Applications*, 273:126833, 2025. ISSN 0957-4174. doi: <https://doi.org/10.1016/j.eswa.2025.126833>. URL <https://www.sciencedirect.com/science/article/pii/S0957417425004555>. [p4]
- M. Nanda, K. Seminar, D. Nandika, and A. Maddu. A comparison study of kernel functions in the support vector machine and a comparison study of kernel functions in the support vector machine and its application for termite detection, 2018. [p3]
- G. Pagallo and D. Haussler. Boolean feature discovery in empirical learning. *Machine Learning*, 5:71–99, 1990. URL <https://api.semanticscholar.org/CorpusID:5661437>. [p1]
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. [p4]
- J. R. Quinlan. Induction of decision trees. *Machine Learning*, 1:81–106, 1986. URL <https://api.semanticscholar.org/CorpusID:189902138>. [p1]
- J. R. Quinlan. *C4. 5: programs for machine learning*. Elsevier, 2014. [p1, 5]
- R. Rivera-Lopez and J. Canul-Reich. Construction of near-optimal axis-parallel decision trees using a differential-evolution-based approach. *IEEE Access*, PP:1–1, 01 2018. doi: 10.1109/ACCESS.2017.2788700. [p1]
- C. Schaffer. Deconstructing the digit recognition problem. In D. Sleeman and P. Edwards, editors, *Machine Learning Proceedings 1992*, pages 394–399. Morgan Kaufmann, San Francisco (CA), 1992. ISBN 978-1-55860-247-2. doi: <https://doi.org/10.1016/B978-1-55860-247-2.50056-5>. URL <https://www.sciencedirect.com/science/article/pii/B9781558602472500565>. [p1]
- W. N. Street, W. H. Wolberg, and O. L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. In *Electronic imaging*, 1993. URL <https://api.semanticscholar.org/CorpusID:14922543>. [p8]
- F. TAKAHASHI and S. Abe. Decision-tree-based multiclass support vector machines. volume 3, pages 1418 – 1422 vol.3, 12 2002. ISBN 981-04-7524-1. doi: 10.1109/ICONIP.2002.1202854. [p4]
- V. Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013. [p3, 4]
- P. C. Young. Recursive estimation and time-series analysis: An introduction. 1984. URL <https://api.semanticscholar.org/CorpusID:60181335>. [p2]

Aneesh Agarwal  
Monash University

[aaga0022@student.monash.edu](mailto:aaga0022@student.monash.edu)

Jack Jewson  
Monash University  
Department of Econometrics and Business Statistics, Monash University, Australia  
[Jack.Jewson@monash.edu](mailto:Jack.Jewson@monash.edu)

*Erik Sverdrup*  
*Monash University*  
*Department of Econometrics and Business Statistics, Monash University, Australia*  
[Erik.Sverdrup@monash.edu](mailto:Erik.Sverdrup@monash.edu)