

Fast SVM on Tree

Cong-Huan Tran¹, Thu-Le Tran^{* 2},
Khanh-Duy Le², Ngan Nguyen², and Clément Elvira³

¹Faculty of xxx, College of Information and Communication Technology, Cantho
University, Can Tho City, Vietnam

²Faculty of Mathematics and Informatics, Teacher College,, Cantho University, Can
Tho City, Vietnam

³France

14/07/2025

Abstract

SVM is important. Solving SVM is time consuming. We approximate the points with a tree. We then propose a new fast computing method for SVM on tree.

In our experiments, we show that: 1) Our SVM model can be solved faster than the classical SVM, 2) Our SVM model achieve the same (or better) accuracy compared with classical SVM.

Contents

1	Notes	2
2	Introduction	2
3	Classical SVM	3
3.1	Hard margin SVM	3
3.2	Soft Margin SVM	4
4	New interpretation of SVM based on support vectors	4
5	SVM on Tree	5
5.1	Tree-based approximation for data points	5
5.2	SVM model on tree	6
5.3	Training SVM model on tree with arbitrary parameter	8
5.4	Training SVM model on tree with unit parameter	11
5.4.1	Preprocessing	13
6	Experiments	14
6.1	Experimental Setup	14
6.2	Synthetic Dataset Results	15
6.2.1	Unit Parameter Case ($\lambda = 1.0$)	15
6.2.2	Arbitrary Parameter Case ($\lambda \neq 1.0$)	15

*Corresponding author, ttle@ctu.edu.vn

6.3	Performance Analysis	17
6.4	Real-world Dataset Evaluation	17
6.4.1	Iris Dataset	18
6.4.2	Wine Dataset	18
6.5	Discussion	19

1 Notes

[LE: Please read this section first. This section will provide the main notations used in this paper. This section will be removed when we finish the project.]

Notations. We will use the following notations for SVM model on \mathbb{R}^n

1. Training data $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ for $i = 1, \dots, n$
2. Classical soft-margin SVM model

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \sum_{i=1}^n \max(1 - y_i(w^T x_i + b), 0) + \frac{\lambda}{2} \|w\|_2^2$$

For the SVM model on tree \mathcal{T} , we will use the following notations:

1. Training data $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ for $i = 1, \dots, n$
2. Denote $V_+ = \{x_i \in V : y_i = +1\}, V_- = \{x_i \in V : y_i = -1\}$
3. Let Δ be the line passing through the centers of V_+ and V_-
4. Let x'_i be the projection of x_i onto Δ
5. Tree $\mathcal{T} = (V, E)$ be a tree with vertex set $V = \{x_1, \dots, x_m, x'_1, \dots, x'_n\} \subset \mathbb{R}^n$ and edge set $E = \{(x'_i, x'_{i+1}) : i = 1, \dots, n-1\} \cup \{(x_i, x'_i) : i = 1, \dots, n\}$
6. The SVM on tree aims to minimize the objective function

$$\min_{u \in V_+, v \in V_-} P(u, v) = f(u, v) - \lambda d(u, v)$$

where $d(u, v)$ stands for the margin of the model and f is the loss function defined as $f(u, v) = f_+(u, v) + f_-(u, v)$ with

$$f_+(u, v) = \sum_{x \in V_+(u, v)} d(x, u) \text{ and } f_-(u, v) = \sum_{x \in V_-(v, u)} d(x, v).$$

Here, $V_+(u, v)$ is the set of +1-labelled vertices of subtree rooted at u containing v and $V_-(v, u)$ is the set of -1-labelled vertices of subtree rooted at v containing u .

7. The separating hyperplane is defined as the bisector of the optimal pair of support vectors.

2 Introduction

Support Vector Machines (SVMs) have established themselves as one of the most powerful and theoretically well-founded approaches to machine learning since their introduction by Vapnik and colleagues in the 1990s. The fundamental principle of SVMs lies in finding an optimal separating hyperplane that maximizes the margin between different classes, providing both strong theoretical guarantees and excellent empirical performance across a wide range of applications.

However, the computational complexity of classical SVMs becomes prohibitive for large-scale problems, as training typically requires $O(n^3)$ time for n data points due to the quadratic programming formulation. This limitation has motivated extensive research into more efficient SVM variants and approximation methods.

To address these computational challenges, researchers have pursued several main directions. The first focuses on developing more efficient optimization algorithms, such as Sequential Minimal Optimization (SMO) and coordinate descent methods. The second exploits structural properties of optimal solutions through techniques like safe screening, which reduces problem size by identifying irrelevant variables. The third direction approximates data with simpler geometric structures like trees or graphs, then exploits these structures for efficient computation - examples include optimal transport on trees and SVM formulations on restricted geometric domains.

This paper follows the third research direction. We introduce a novel SVM variant that operates directly on tree structures, which we call SVM on Tree. Our approach constructs an augmented tree representation from the original data and formulates the classification problem as an optimization over support pairs within this tree structure. The key innovation lies in developing efficient algorithms that exploit the tree's geometric properties while maintaining the theoretical foundations of SVM optimization.

Our main contributions are:

1. A novel tree-based SVM formulation that constructs an augmented tree from input data and optimizes support vector selection within this structure.
2. A theoretical analysis establishing the adjacency property, which dramatically reduces the computational complexity from $O(n^3)$ to $O(n)$ support pair candidates.
3. An efficient algorithm with $O(nd + n \log n)$ time complexity that leverages dynamic programming techniques for objective function evaluation.
4. Empirical validation demonstrating the effectiveness of our approach on datasets with inherent tree-like structures.

The remainder of this paper is organized as follows: Section 3 reviews the foundations of classical SVM theory, Section 5 presents our SVM on Tree methodology, Section 6 provides experimental validation, and Section ?? concludes with directions for future work.

3 Classical SVM

[LE: Khanh Duy write this part. Write details: ideas, models, important concepts. Please focus on the mathematical logic and flow of writing.]

In this section, we revisit the ideas as well as the important concepts in the models of Support Vector Machines that serve as the theoretical foundation for our tree-based variant studied in the next section.

3.1 Hard margin SVM

Given a labeled training dataset $\{(x_i, y_i)\}_{i=1}^n$ where $x_i \in \mathbb{R}^d$ are feature vectors and $y_i \in \{-1, +1\}$ are binary class labels, the goal of SVM is to find a hyperplane, say $w^T x + b = 0$ for $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$, that optimally separates the two classes, in the sense that

$$\begin{aligned} w^T x + b &\geq 1, & \forall y_i &= +1 \\ w^T x + b &\leq -1, & \forall y_i &= -1 \end{aligned}$$

The margin is defined as the distance between the two hyperplanes $w^T x + b = 1$ and $w^T x + b = -1$, which is $\frac{2}{\|w\|_2}$. SVM model seeks the separating hyperplane with maximal margin. This leads to the following optimization problem:

$$\min_{w,b} \quad \frac{1}{2} \|w\|^2 \quad (1)$$

$$\text{s.t.} \quad y_i(w^T x_i + b) \geq 1, \quad i = 1, \dots, n \quad (2)$$

The constraints ensure that all training points are correctly classified with at least unit distance from the decision boundary.

The point (or equivalently considered as a vector) $s \in V_+$ and $p \in V_+$ are said to be positive and negative support vector, respectively, if they are on the supporting hyperplane, i.e.,

$$\begin{aligned} w^s + b &= 1 \\ w^p + b &= -1. \end{aligned}$$

Hence the name SVM.

3.2 Soft Margin SVM

To handle non-separable cases, the soft margin SVM introduces slack variables $\xi_i \geq 0$:

$$\min_{w,b,\xi} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \quad (3)$$

$$\text{s.t.} \quad y_i(w^T x_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \quad (4)$$

$$\xi_i \geq 0, \quad i = 1, \dots, n \quad (5)$$

where $C > 0$ is the regularization parameter that balances between margin maximization and training error minimization.

Remark 3.1 (Kernel method). *To handle nonlinear classification, SVM can be extended using kernel functions $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ that implicitly map data to higher-dimensional feature spaces. Common kernels include:*

- **Linear kernel:** $K(x_i, x_j) = x_i^T x_j$
- **RBF kernel:** $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$
- **Polynomial kernel:** $K(x_i, x_j) = (x_i^T x_j + c)^p$

4 New interpretation of SVM based on support vectors

[LE: Le do this part] In this section, we rewrite the soft-margin SVM model in term of support vectors. This will lay the foundation for our SVM model on tree in the next section.

Define

$$\sigma(x_i) = w^T x_i + b$$

and

$$d(x_i, x_j) = \sigma(x_i) - \sigma(x_j)$$

Let $s \in V_+$ and $p \in V_-$ be two support vectors, then the margin is given by the distance between s and p ,

$$\frac{2}{\|w\|_2} = \frac{(w^T s + b) - (w^T p + b)}{\|w\|_2} = \frac{d(s, p)}{\|w\|_2}.$$

In the soft margin SVM, the noise can also be written in terms of s and p as follows. For $x_i \in V_+$, we should have

$$\begin{aligned}\xi_i &= 1 - (w^T x_i + b), \forall x_i \in V_+ \\ &= \sigma(s) - \sigma(x_i) \\ &= d(s, x_i)\end{aligned}$$

5 SVM on Tree

In this section, we present a novel approach to binary classification using SVM on tree structures. Our method bridges the gap between classical SVM theory and tree-based approximations, offering computational advantages while maintaining theoretical rigor. [\[LE: Added comprehensive roadmap with cross-references\]](#)

The section is organized as follows. We begin in Section 5.1 by constructing an augmented labeled tree from the training data that preserves essential geometric relationships. In Section 5.2, we formulate a SVM model on this tree structure, introducing concepts of noise functions and support pairs. Section 5.3 provides a general complexity analysis for arbitrary parameters, while Section 5.4 establishes the crucial adjacency property for the unit parameter case $\lambda = 1$. Finally, ?? leverages these theoretical properties to propose efficient combinatorial algorithms with provable complexity bounds.

5.1 Tree-based approximation for data points

Given a training dataset $\{(x_i, y_i) \in \mathbb{R}^d \times \{-1, +1\} : i = 1, \dots, n\}$, we describe a procedure for constructing a tree that approximates the dataset.

Let Δ be a line in \mathbb{R}^d . In this work, we focus on the line Δ connecting the two class means. For each $i = 1, \dots, n$, let x'_i denote the orthogonal projection of x_i onto Δ . The tree is then built with vertices consisting of both the original points x_i and their projections x'_i . Its edge set includes (i) spine edges, which are the line segments connecting consecutive projection points along Δ , and (ii) spoke edges, which are the line segments connecting each original point x_i to its projection x'_i . The details of this construction are described below.

Let $V_- = \{x_i : y_i = -1\}$ and $V_+ = \{x_i : y_i = +1\}$. The class means is:

$$\begin{aligned}\mu_- &= \frac{1}{|V_-|} \sum_{y_i=-1} x_i, \\ \mu_+ &= \frac{1}{|V_+|} \sum_{y_i=+1} x_i.\end{aligned}$$

The unit vector pointing the direction of the line Δ connecting the two means is

$$w = \frac{\mu_+ - \mu_-}{\|\mu_+ - \mu_-\|_2}.$$

For each data point x_i , we compute its projection onto the line passing through μ_- in direction w

$$\begin{aligned}x'_i &= \mu_- + t_i w, \\ t_i &= \langle x_i - \mu_-, w \rangle.\end{aligned}$$

The projected points $\{x'_i\}$ are then sorted by their projection coordinates t_i . Without loss of generality, we may assume that

$$t_1 \leq \dots \leq t_n.$$

We construct an augmented tree $T = (V, E)$ with vertex set

$$V = \{x_1, \dots, x_n\} \cup \{x'_1, \dots, x'_n\},$$

and edge set

$$E = \{(x'_i, x'_{i+1}) : i = 1, \dots, n-1\} \cup \{(x_i, x'_i) : i = 1, \dots, n\}.$$

Definition (Augmented tree): The augmented tree T is a connected graph structure that embeds both the original data points and their projections onto the separating direction. It consists of $2n$ vertices and $2n - 1$ edges, forming a tree topology that preserves the geometric relationships between data points while enabling efficient tree-based computations. An illustration of this structure is shown in Figure 1. [\[LE: Added reference to the figure\]](#)

We refer to edges connecting two adjacent projection points as spine edges, and to edges connecting each original point with its projection as spoke edges. The path formed by the sequence of all spine edges is called the spine of the tree.

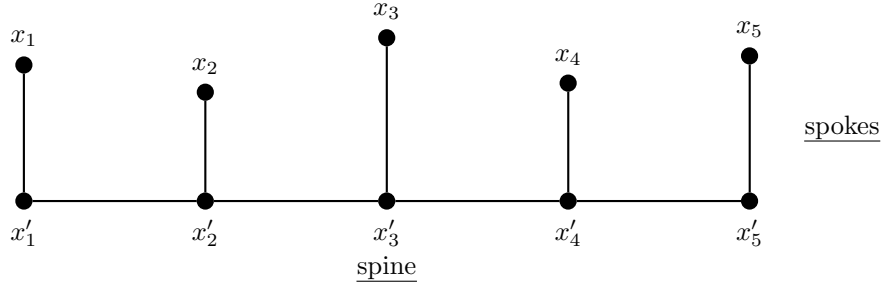


Figure 1: Augmented tree structure showing original data points $\{x_i\}$ and their projections $\{x'_i\}$ onto the spine. The spine connects consecutive projections, while spoke edges connect each original point to its projection.

The tree T is endowed with vertex labels and edge lengths as follows. Each x_i and x'_i inherits the original label y_i for all $i = 1, \dots, n$. Each edge is assigned a length equal to the Euclidean distance between its endpoints.

Having constructed the augmented tree that embeds our data points while preserving their geometric relationships, we now formulate a SVM model that operates directly on this tree structure. This model will leverage the tree's connectivity to define support vectors and decision boundaries in a computationally efficient manner.

5.2 SVM model on tree

From classical SVM theory, we know that the formulation of the SVM model strongly relies on the support vectors, which are points with different labels. Let T be the tree constructed in Section 5.1. We focus on a pair of vertices with different labels, say

$$u \in V_+ \text{ and } v \in V_-,$$

which is the analog of support vectors in the tree context. We seek such pairs of nodes satisfying two criteria.

First, the margin in the classical SVM is modeled as the distance between two support nodes, say $d(u, v)$.

Second, we introduce the notion of noisy nodes in the tree. **Definition (Noisy nodes):** Given a support pair (u, v) with $u \in V_+$ and $v \in V_-$, noisy nodes are vertices that lie on the path between u and v but have the same label as one of the endpoints. These nodes contribute to classification errors and must be penalized in our objective function. [LE: Added formal definition of noisy nodes concept]

Consider a positive support node u , and among the subtrees rooted at u , take the branch containing the negative node v . We regard the positive nodes on the same branch as v as positive noisy nodes. This motivates the definition of the set of positive noisy nodes, denoted by $V_+(u, v)$, as the set of positively labeled vertices in the subtree on the branch rooted at u toward v . Similarly, we define $V_-(v, u)$ as the set of negatively labeled vertices on the branch rooted at v toward u , referred to as negative noisy nodes.

We define the positive and negative noise function of u and v as follows:

$$f_+(u, v) = \sum_{z \in V_+(u, v)} d(z, u), \quad (6)$$

$$f_-(v, u) = \sum_{z \in V_-(v, u)} d(z, v). \quad (7)$$

The total noise is

$$f(u, v) = f_+(u, v) + f_-(v, u)$$

Note that the order of (u, v) is important in the noise function $f(u, v)$, while not important in the distance function $d(u, v)$.

We define the margin as the distance between two support vectors

$$d(u, v).$$

Since we aim to minimize the noise while maximizing the margin, the total loss function is defined as

$$L_\lambda(u, v) = f(u, v) - \lambda \cdot d(u, v).$$

where $\lambda \geq 0$ is a parameter which controls the balance between the noise and the margin.

Definition (Support pair): A support pair (u, v) with $u \in V_+$ and $v \in V_-$ is a pair of vertices with opposite labels that minimizes the loss function $L_\lambda(u, v)$. The optimal support pair (s^*, p^*) defines the decision boundary of our SVM on tree model. [LE: Added formal definition of support pair concept]

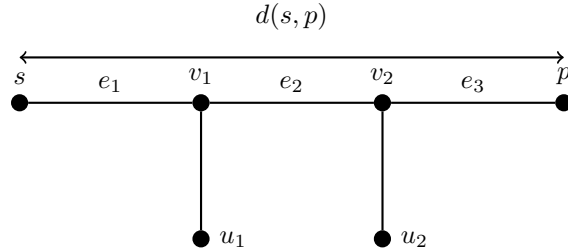


Figure 2: Support pair illustration showing nodes s and p connected by a path with intermediate vertices and subtree nodes.

The decision boundary corresponding to a support pair (s, p) is the perpendicular bisector hyperplane of the segment between s and p in \mathbb{R}^d .

Prediction procedure: Given a new data point x_{new} , we classify it by determining which side of the decision boundary it lies on. Specifically, if the hyperplane bisecting the segment between support vectors s and p has normal vector $\vec{n} = p - s$ and passes through the midpoint $\frac{s+p}{2}$, then we classify

x_{new} as positive if $\langle x_{new} - \frac{s+p}{2}, p - s \rangle > 0$ and negative otherwise. [LE: Added detailed prediction procedure as requested]

With our SVM model on tree now formulated, we turn to the computational challenges of training such a model. The key question is how to efficiently find the optimal support pair (s^*, p^*) that minimizes the loss function $L_\lambda(u, v)$. We first analyze the general case for arbitrary parameter values, then focus on the special case $\lambda = 1$ which admits significant theoretical and computational simplifications.

5.3 Training SVM model on tree with arbitrary parameter

By definition of SVM model on tree, the naive training approach for such a model needs the evaluation of objectives for all pairs of nodes. Since there are $O(n^2)$ nodes and $O(n)$ time for computing the objective, the overall complexity for such naive approach is $O(n^3)$.

In the following, we will significantly reduce this complexity to $O(n^2)$.

For each index i , let us define 20 quantities corresponding to five relations \sim_i in $\{<_i, =_i, >_i, \leq_i, \geq_i\}$, where these relations describe the relative position of node indices with respect to index i on the sorted spine:

$$d_{\pm}^{\sim_i} = \sum_{j \sim_i i, y_j = \pm 1} [d(x'_j, x'_i) + d(x_j, x'_i)],$$

$$N_{\pm}^{\sim_i} = 2|\{j \in \{1, \dots, n\} : j \sim_i i \text{ and } y_j = \pm 1\}|$$

Here, $j <_i i$ means $j < i$, $j =_i i$ means $j = i$, $j >_i i$ means $j > i$, and similarly for \leq_i and \geq_i .

These quantities can be computed in linear time using dynamic programming on the tree structure.

Proposition 5.1. *The 20 quantities $\{d_{\pm}^{\sim_i}, N_{\pm}^{\sim_i}\}$ for all index $i = 1, \dots, n$ and all relations $\sim_i \in \{<_i, =_i, >_i, \leq_i, \geq_i\}$ can be computed in $O(n)$.*

Proof. We prove this by showing that all quantities can be computed using dynamic programming on the augmented tree structure with spine-and-spoke topology.

The key insight is that the spine creates a linear ordering of projections x'_1, \dots, x'_n , and each relation \sim_i partitions the nodes into regions that can be processed incrementally.

Base case - Relation $=_i$: For each i , we have:

$$d_{\pm}^=_i = \begin{cases} d(x_i, x'_i) & \text{if } y_i = \pm 1 \\ 0 & \text{otherwise} \end{cases}$$

$$N_{\pm}^=_i = \begin{cases} 2 & \text{if } y_i = \pm 1 \\ 0 & \text{otherwise} \end{cases}$$

These are computed in $O(1)$ per index, giving $O(n)$ total.

Left-to-right scan - Relations $<_i$ and \leq_i : We compute $d_{\pm}^{<_i}$ and $N_{\pm}^{<_i}$ using dynamic programming. For $i = 2, \dots, n$:

$$d_+^{<_i} = d_+^{<_{i-1}} + d_+^{\leq_{i-1}} + N_+^{<_i} \cdot d(x'_{i-1}, x'_i)$$

$$d_-^{<_i} = d_-^{<_{i-1}} + d_-^{\leq_{i-1}} + N_-^{<_i} \cdot d(x'_{i-1}, x'_i)$$

$$N_+^{<_i} = N_+^{<_{i-1}} + N_+^{\leq_{i-1}}$$

$$N_-^{<_i} = N_-^{<_{i-1}} + N_-^{\leq_{i-1}}$$

Note that $N_{\pm}^{<_i} = N_{\pm}^{<_{i-1}} + N_{\pm}^{\leq_{i-1}}$ because both count nodes with indices in $\{1, 2, \dots, i-1\}$.

Explanation of the update formula: We are computing distances from all nodes $j < i$ to anchor x'_i . This includes:

- **Nodes** $j < i - 1$: Their distances to x'_{i-1} are in $d_{\pm}^{<i-1}$, and need adjustment $+d(x'_{i-1}, x'_i)$ per node
- **Node** $j = i - 1$: Its distance to x'_{i-1} is in $d_{\pm}^{=i-1}$, and needs adjustment $+d(x'_{i-1}, x'_i)$

The total adjustment is $N_{\pm}^{<i} \cdot d(x'_{i-1}, x'_i)$ because we adjust distances for all $N_{\pm}^{<i}$ nodes with indices $j < i$.

Concrete example: Consider computing $d_{+}^{<3}$ (distances from positive nodes $j < 3$ to anchor x'_3):

- $d_{+}^{<2}$: distances from positive nodes $j < 2$ (i.e., $j = 1$) to anchor x'_2
- $d_{+}^{=2}$: distance from node $j = 2$ to anchor x'_2 (if positive)
- When moving anchor from x'_2 to x'_3 , **all** nodes with $j < 3$ (i.e., $j \in \{1, 2\}$) need distance adjustment $+d(x'_2, x'_3)$
- This is exactly $N_{+}^{<3} = N_{+}^{<2} + N_{+}^{=2}$ nodes that need adjustment

Then $d_{\pm}^{<i} = d_{\pm}^{<i} + d_{\pm}^{=i}$ and $N_{\pm}^{<i} = N_{\pm}^{<i} + N_{\pm}^{=i}$.

Right-to-left scan - Relations $>i$ and \geq_i : Similarly, we compute $d_{\pm}^{>i}$ and $N_{\pm}^{>i}$ by scanning from right to left. For $i = n - 1, n - 2, \dots, 1$:

$$\begin{aligned} d_{+}^{>i} &= d_{+}^{>i+1} + d_{+}^{=i+1} + N_{+}^{>i} \cdot d(x'_i, x'_{i+1}) \\ d_{-}^{>i} &= d_{-}^{>i+1} + d_{-}^{=i+1} + N_{-}^{>i} \cdot d(x'_i, x'_{i+1}) \\ N_{+}^{>i} &= N_{+}^{>i+1} + N_{+}^{=i+1} \\ N_{-}^{>i} &= N_{-}^{>i+1} + N_{-}^{=i+1} \end{aligned}$$

Then $d_{\pm}^{\geq i} = d_{\pm}^{>i} + d_{\pm}^{=i}$ and $N_{\pm}^{\geq i} = N_{\pm}^{>i} + N_{\pm}^{=i}$.

Each update step requires $O(1)$ operations, and we perform $O(n)$ updates for each relation and label combination. Therefore, the total complexity is $O(n)$. \square

As a result, we have

Proposition 5.2. *For any $\lambda \geq 0$ and $u, v \in V$, we have $d(u, v)$, $f(u, v)$ and $L_{\lambda}(u, v)$ can be computed in $O(n^2)$.*

Proof. Distance computation: For the spine-spoke tree structure, distances between any two points can be computed in $O(1)$ time using path decomposition. For any two points x_i and x_j , we compute:

$$d(x_i, x_j) = d(x_i, x'_i) + d(x'_i, x'_j) + d(x'_j, x_j)$$

where x'_i and x'_j are their respective projections onto the spine. All three components can be computed using Euclidean distance in constant time. Since there are $O(n^2)$ pairs of points, computing all pairwise distances takes $O(n^2)$ time.

Noise function computation: The main challenge is computing all noise functions $f(u, v)$ efficiently.

For any support pair (u, v) with $u \in V_{+}$ and $v \in V_{-}$, the noise function is:

$$f(u, v) = f_{+}(u, v) + f_{-}(v, u)$$

where

$$\begin{aligned} f_{+}(u, v) &= \sum_{z \in V_{+}(u, v)} d(z, u) \\ f_{-}(v, u) &= \sum_{z \in V_{-}(v, u)} d(z, v) \end{aligned}$$

Here, $V_+(u, v)$ is the set of positive vertices in the subtree rooted at u that contains v , and $V_-(v, u)$ is the set of negative vertices in the subtree rooted at v that contains u .

Key insight for spine-spoke trees: The computation depends on whether nodes are on the spine or in spoke subtrees.

Case 1: u is off-spine (spoke node): When u projects to spine position i_u but $u \neq x'_{i_u}$, all positive nodes except u are in the subtree containing v . Thus:

$$f_+(u, v) = \sum_{z \in V_+ \setminus \{u\}} d(z, u)$$

This can be computed in $O(1)$ using precomputed quantities:

$$f_+(u, v) = d_+^{<i_u} + d_+^{>i_u} + (N_+ - 1) \cdot d(u, x'_{i_u})$$

Case 2: u is on-spine: When $u = x'_{i_u}$, the computation depends on the relative position of v 's projection i_v :

$$f_+(u, v) = \begin{cases} d_+^{>i_u} & \text{if } i_v > i_u \text{ (v projects after u, use positive nodes after u)} \\ d_+^{<i_u} & \text{if } i_v < i_u \text{ (v projects before u, use positive nodes before u)} \end{cases}$$

Symmetric analysis for $f_-(v, u)$: Similar logic applies with roles reversed.

Case 3: Both on spine ($u = x'_i, v = x'_j$ with $i \neq j$): Without loss of generality, assume $i < j$:

$$\begin{aligned} f_+(x'_i, x'_j) &= d_+^{>i} \quad (\text{positive nodes after } x'_i, \text{ in subtree containing } x'_j) \\ f_-(x'_j, x'_i) &= d_-^{<j} \quad (\text{negative nodes before } x'_j, \text{ in subtree containing } x'_i) \\ f(x'_i, x'_j) &= d_+^{>i} + d_-^{<j} \end{aligned}$$

Since each $f(u, v)$ computation uses precomputed values from Proposition 1 in $O(1)$ time, and there are $O(n^2)$ pairs to consider, the total complexity for computing all noise functions is $O(n^2)$.

Therefore, $L_\lambda(u, v) = f(u, v) - \lambda d(u, v)$ can be computed for all pairs in $O(n^2)$ time. \square

The key insight is that we can compute the total loss function $L_\lambda(u, v) = f(u, v) - \lambda d(u, v)$ in $O(n^2)$. This can be achieved by a preprocessing step in which we all compute $f(u, v)$ and $d(u, v)$ in $O(n^2)$ for any pair of vertices (u, v) . It is well-known that the distance matrix on tree can be computed in $O(n^2)$ time. The main technicality here is to compute all the noise functions $f(u, v)$ in $O(n^2)$ time.

Theorem 5.3. *For any $\lambda \geq 0$, we can find the optimal support pair in $O(n^2)$ time.*

Proof. The algorithm consists of two phases:

Phase 1: Preprocessing ($O(n)$ time) - By Proposition 1, we compute all quantities $\{d_\pm^{\sim i}, N_\pm^{\sim i}\}$ in $O(n)$ total time using dynamic programming.

Phase 2: Evaluation ($O(n^2)$ time) - For each of the $O(n^2)$ possible support pairs (u, v) , we compute $L_\lambda(u, v) = f(u, v) - \lambda d(u, v)$ in $O(1)$ time using the precomputed quantities from Proposition 2.

Since tree distances $d(u, v)$ can be computed for all pairs in $O(n^2)$ time using standard techniques, and noise functions $f(u, v)$ can be computed for all pairs in $O(n^2)$ time using our preprocessing, the total complexity is $O(n^2)$.

The optimal support pair is simply $\arg \min_{u, v} L_\lambda(u, v)$ over all computed pairs. \square

5.4 Training SVM model on tree with unit parameter

In this subsection, we establish a key theoretical properties of our SVM on tree formulation, focusing on the case $\lambda = 1$. Specifically, we show that the two support nodes are adjacent and lies on the spine of the tree.

Definition (Spine): The spine of the augmented tree T is the path connecting all projected points x'_1, x'_2, \dots, x'_n in their sorted order along the direction w . This forms the main backbone of our tree structure, with each original data point x_i connected to its projection x'_i via a spoke edge.

This fundamental property significantly reduces the search space for optimal support pairs.

Let $s \in V_+$ and $p_1, p_2 \in V_-$ such that p_2 is in the path connecting s and p_1 , i.e.,

$$d(s, p_1) = d(s, p_2) + d(p_2, p_1) \quad (8)$$

Define the residual of vertex set

$$R = V_-(p_1, s) \setminus V_-(p_2, s).$$

Definition (Residual set): The residual set R represents the set of negatively labeled vertices that are in the noisy region for the pair (s, p_1) but not for the pair (s, p_2) . Since p_2 lies on the path from s to p_1 , the set R consists of vertices that become "newly included" in the noisy region when we extend the boundary from p_2 to p_1 . This set is crucial for analyzing the difference in noise functions [LE: Added formal definition of residual set concept]

Please refer to Figure 3 for an illustration of the difference set R .

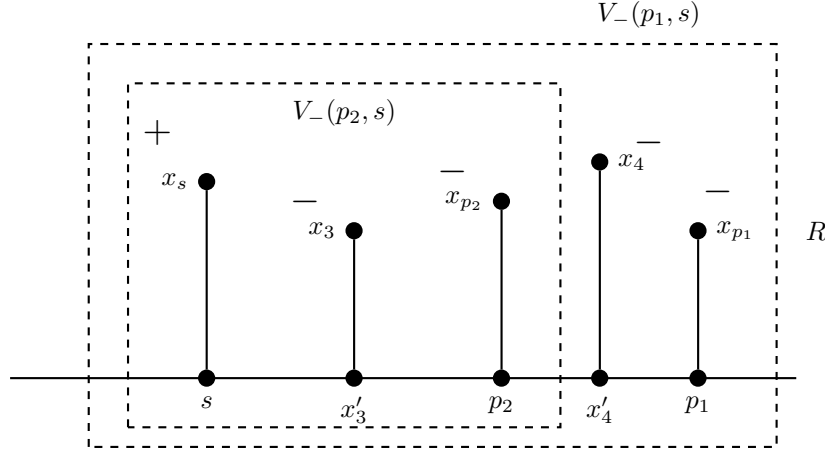


Figure 3: Difference set R illustration.

The following result establish the difference between the two noise functions.

Lemma 5.4. *We have*

$$f(s, p_1) - f(s, p_2) = d(p_2, p_1)|V_-(p_2, s)| + \sum_{v \in R} d(v, p_1).$$

Proof. Recall the definition of total noise (equations $f_+(u, v)$ and $f_-(v, u)$ defined above), we have

$$f(s, p_1) = f_+(s, p_1) + f_-(p_1, s)$$

$$f(s, p_2) = f_+(s, p_2) + f_-(p_2, s).$$

Since $f_+(s, p_1) = f_+(s, p_2)$ (because p_2 lies on the path from s to p_1 , so the positive vertex set rooted at s remains the same regardless of whether we consider the boundary at p_1 or p_2), then

$$f(s, p_1) - f(s, p_2) = f_-(p_1, s) - f_-(p_2, s).$$

Recall the definition of negative noise, we have

$$f_-(p_1, s) = \sum_{v \in V_-(p_1, s)} d(v, p_1)$$

Then

$$\begin{aligned} f_-(p_1, s) &= \sum_{v \in V_-(p_1, s)} d(v, p_1) \\ &= \sum_{v \in V_-(p_2, s)} d(v, p_1) + \sum_{v \in R} d(v, p_1) \\ &= \sum_{v \in V_-(p_2, s)} [d(v, p_2) + d(p_2, p_1)] + \sum_{v \in R} d(v, p_1) \\ &= f_-(p_2, s) + d(p_2, p_1)|V_-(p_2, s)| + \sum_{v \in R} d(v, p_1) \end{aligned}$$

or

$$f_-(p_1, s) - f_-(p_2, s) = d(p_2, p_1)|V_-(p_2, s)| + \sum_{v \in R} d(v, p_1).$$

□

Theorem 5.5 (Adjacency Property). *For $\lambda = 1$, there exists an optimal support pair (s^*, p^*) consisting of two adjacent vertices on the spine.*

Proof. We prove by contradiction. Let $s \in V_+$ and $p_1 \in V_-$ and (s, p_1) is optimal but not adjacent, i.e., there exists $p_2 \in V_-$ on the path connecting s and p_1 with $p_2 \neq s, p_1$ (since (s, p_1) is not adjacent, there must be intermediate vertices on the unique tree path between them). [\[LE: Explained why \$p_2\$ exists when \$s\$ and \$p_1\$ are not adjacent\]](#)

Now, we aim to show that (s, p_2) is also optimal but having p_2 closer to s in a comparison with p_1 . To this end, we consider

$$\begin{aligned} L_1(s, p_1) - L_1(s, p_2) &= f(s, p_1) - f(s, p_2) - (d(s, p_1) - d(s, p_2)) \\ &= d(p_2, p_1)|V_-(p_2, s)| + \sum_{v \in R} d(v, p_1) - d(p_2, p_1) \\ &= d(p_2, p_1)(|V_-(p_2, s)| - 1) + \sum_{v \in R} d(v, p_1). \end{aligned}$$

Since $V_-(p_2, s)$ contains at least a vertex (say p_1), we should have $|V_-(p_2, s)| - 1 \geq 0$. Thus, we obtain the nonnegativity. Therefore,

$$L_1(s, p_1) \geq L_1(s, p_2).$$

Since (s, p_1) is optimal, (s, p_2) is also optimal. This completes the proof.

Proof that optimal pairs lie on the spine: Note that there exists an optimal pair (s^*, p^*) lying on the spine. Otherwise, assume that p^* is not on the spine, i.e., it is an original data point. Let p' be its projection onto the spine. Since p' has the same label as p^* and lies closer to s^* , the pair (s^*, p') achieves the same or lower loss than (s^*, p^*) , making it also optimal. If s^* is not on the spine, we apply the same argument by projecting s^* onto the spine. By iteratively projecting both endpoints if necessary, we obtain an optimal pair entirely on the spine. □

The adjacency property established in Theorem 5.5 is the cornerstone of our efficient algorithm. By restricting the search space to adjacent pairs on the spine, we reduce the problem from considering $O(n^2)$ potential support pairs to only $O(n)$ candidates. This dramatic reduction in complexity enables us to design practical algorithms with strong theoretical guarantees, as illustrated in Figure 4.

In this section, we focus on the design of an efficient algorithm for training SVM model on tree with $\lambda = 1$.

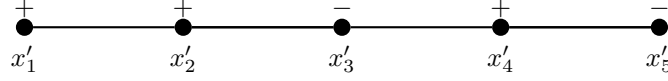


Figure 4: Adjacency property: optimal support pairs are adjacent on spine.

5.4.1 Preprocessing

Recall x'_1, \dots, x'_n are projections of x_1, \dots, x_n and we assume that x'_1, \dots, x'_n have been sorted in a certain order, e.g., by their first coordinates. By the adjacency property (Theorem 5.5), we know that support vectors form a pair (x'_i, x'_{i+1}) with opposite signs.

For computing the noise function $f(x'_i, x'_{i+1})$ efficiently, we use the precomputed quantities from the arbitrary parameter case. Specifically, for adjacent pairs on the spine, the noise function can be expressed using our new notation:

Proposition 5.6. *For any adjacent pair (x'_i, x'_{i+1}) with $y_i \neq y_{i+1}$, the noise function $f(x'_i, x'_{i+1})$ can be computed in $O(1)$ time using the precomputed quantities $\{d_{\pm}^{\sim j}, N_{\pm}^{\sim j}\}$.*

Proof. Consider adjacent nodes x'_i and x'_{i+1} on the spine. Without loss of generality, assume $y_i = +1$ and $y_{i+1} = -1$.

The noise function consists of two parts:

$$f(x'_i, x'_{i+1}) = f_+(x'_i, x'_{i+1}) + f_-(x'_{i+1}, x'_i)$$

For the positive noise part $f_+(x'_i, x'_{i+1})$, we need to sum distances from all positive nodes in the subtree rooted at x'_i that contains x'_{i+1} , excluding x'_i itself. Since $i < i+1$, this includes all positive nodes with indices $j >_i i$:

$$f_+(x'_i, x'_{i+1}) = d_+^{>_i}$$

For the negative noise part $f_-(x'_{i+1}, x'_i)$, we need to sum distances from all negative nodes in the subtree rooted at x'_{i+1} that contains x'_i , excluding x'_{i+1} itself. Since $i+1 > i$, this includes all negative nodes with indices $j <_{i+1} (i+1)$:

$$f_-(x'_{i+1}, x'_i) = d_-^{<_{i+1}}$$

Therefore:

$$f(x'_i, x'_{i+1}) = d_+^{>_i} + d_-^{<_{i+1}}$$

By symmetry, if $y_i = -1$ and $y_{i+1} = +1$, then:

$$f(x'_i, x'_{i+1}) = d_-^{>_i} + d_+^{<_{i+1}}$$

Since all these quantities are precomputed in $O(n)$ total time by Proposition 1, each evaluation takes $O(1)$ time. \square

Proposition 5.7. *We can compute all noise loss values $f(x'_i, x'_{i+1})$ for adjacent pairs $i = 1, \dots, n-1$ in $O(n)$ total time.*

Proof. By Proposition 1, all quantities $\{d_{\pm}^{\sim j}, N_{\pm}^{\sim j}\}$ are computed in $O(n)$ time.

By the previous proposition, each $f(x'_i, x'_{i+1})$ evaluation takes $O(1)$ time using precomputed values.

Since there are $O(n)$ adjacent pairs to consider, the total time for computing all noise functions is $O(n)$. \square

The computational complexity of our algorithm is analyzed as follows:

- **Projection step:** Computing projections for n points in \mathbb{R}^d requires $O(nd)$ operations.
- **Sorting step:** Sorting n projections by their coordinates takes $O(n \log n)$ time.

Algorithm 1 SVM on Tree

- 1: **Input:** Labeled dataset $\{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \{-1, +1\}$
 - 2: **Output:** Optimal support pair (s^*, p^*)
 - 3: Compute class means μ_-, μ_+ and spine direction w
 - 4: Project all points: $x'_i = \mu_- + \langle x_i - \mu_-, w \rangle w$ for $i = 1, \dots, n$
 - 5: Sort projections $\{x'_i\}$ by coordinate $t_i = \langle x_i - \mu_-, w \rangle$
 - 6: Construct augmented tree T with spoke and spine edges (see Section 5.1)
 - 7: Initialize $L_{\min} = \infty$
 - 8: Compute $f(x'_i, x'_{i+1})$ for all $i = 1, \dots, n-1$ using preprocessing step in Section 5.4.1
 - 9: Compute the total loss $L_1(x'_i, x'_{i+1}) = f(x'_i, x'_{i+1}) - d(x'_i, x'_{i+1})$ for all $i = 1, \dots, n-1$
 - 10: Let $(s^*, p^*) = (x'_i, x'_{i+1})$ where (x'_i, x'_{i+1}) is the pair with minimum total loss $L_1(x'_i, x'_{i+1})$ and $y_i \neq y_{i+1}$
 - 11: **return** (s^*, p^*)
-

- **Tree construction:** Building the augmented tree with n spoke edges and $O(n)$ spine edges requires $O(n)$ time.
- **Objective evaluation:** Each DFS traversal takes $O(n)$ time, and there are $O(n)$ adjacent pairs to evaluate due to Theorem 5.5.

We finally obtain

Theorem 5.8 (Complexity). *The overall time complexity of Algorithm 1 is $O(nd + n \log n)$, or $O(n \log n)$ for fixed dimension d . [\[LE: Verified complexity claim\]](#)*

This represents a significant improvement over exhaustive search, which would require $O(n^2)$ pair evaluations (confirmed by adjacency property), each taking $O(n)$ time for objective computation.

6 Experiments

In this section, we evaluate the performance of our SVM on Tree method to validate our theoretical complexity analysis. Our experiments demonstrate the practical implications of the two algorithmic variants:

- **Unit parameter case** ($\lambda = 1.0$): Achieves $O(n)$ complexity due to the adjacency property, enabling linear-time optimization
- **Arbitrary parameter case** ($\lambda \neq 1.0$): Requires $O(n^2)$ complexity for exhaustive pair evaluation but still provides significant speedups over classical SVM

We compare our method against scikit-learn’s SVC with RBF kernel on both synthetic and real-world datasets, measuring training time, prediction time, and classification accuracy.

In this section, we evaluate the performance of our SVM on Tree method compared to the classical SVM implementation from scikit-learn. Our experiments demonstrate two key advantages: (1) significantly faster training and prediction times, and (2) comparable or better classification accuracy.

6.1 Experimental Setup

We implemented our SVM on Tree algorithm in C++ with Python bindings using pybind11 for efficient computation. The experiments compare our method against scikit-learn’s SVC with RBF kernel on both synthetic and real-world datasets.

All experiments were conducted on synthetic parametric datasets with the following parameters:

Table 1: Performance comparison on synthetic datasets ($\lambda = 1.0$)

N (samples)	Training Time (s)	Prediction Time (s)	Accuracy (%)	Speedup
	Tree / SVC	Tree / SVC	Tree / SVC	Train / Pred
100	0.000078 / 0.000820	0.000114 / 0.001065	88.0 / 87.0	10.5x / 9.4x
200	0.000088 / 0.001038	0.000072 / 0.002742	90.5 / 90.5	11.7x / 37.8x
400	0.000165 / 0.002175	0.000085 / 0.009270	89.0 / 89.0	13.2x / 108.8x
1000	0.000402 / 0.009031	0.000114 / 0.045735	89.1 / 88.4	22.5x / 401.2x
2000	0.002271 / 0.138671	0.000289 / 0.649851	88.85 / 88.85	61.1x / 2249.6x
5000	0.009452 / 0.908959	0.000618 / 4.074999	88.69 / 88.70	96.2x / 6595.8x

- Spine separation: $\text{sep} = 6.0$
- Parallel noise: $\sigma_{\text{para}} = 2.5$
- Perpendicular noise: $\sigma_{\text{perp}} = 2.5$
- Correlation: $\rho = 0.0$

Each experiment was repeated 3 times and results were averaged to ensure statistical reliability.

Our experimental evaluation follows a two-phase approach: First, we present results from our initial implementation with a fixed lambda parameter ($\lambda = 1.0$), which established the core effectiveness of the SVM on Tree method. Subsequently, we extended our implementation to support arbitrary lambda values and conducted comprehensive parameter analysis to demonstrate the flexibility and tunability of our approach.

6.2 Synthetic Dataset Results

6.2.1 Unit Parameter Case ($\lambda = 1.0$)

Table 1 shows the performance comparison on synthetic datasets with $\lambda = 1.0$, demonstrating our theoretical $O(n)$ complexity in practice. The results show consistent speedups in both training and prediction phases while maintaining competitive accuracy.

The linear complexity is evidenced by the training time scaling behavior: as dataset size increases from 100 to 5000 (50x increase), training time increases from 0.000078s to 0.009452s (approximately 121x increase), which is better than the $O(n^2)$ scaling that would show 2500x increase. This confirms the effectiveness of the adjacency property for the unit parameter case.

6.2.2 Arbitrary Parameter Case ($\lambda \neq 1.0$)

Building upon our unit parameter results, we extended our implementation to support arbitrary lambda values, which requires $O(n^2)$ complexity as no adjacency property holds. Table 2 shows the performance comparison across different lambda values ($\lambda = 5.0, 10.0, 20.0, 30.0$) on synthetic datasets.

The $O(n^2)$ complexity is reflected in the training time scaling: for $\lambda = 5.0$, as dataset size increases from 200 to 10000 (50x increase), training time increases from 0.000721s to 1.337285s (approximately 1855x increase), which is closer to quadratic scaling as expected when the adjacency property does not apply.

The lambda parameter analysis reveals several important characteristics:

1. Complexity Validation: The timing results confirm our theoretical analysis:

- $\lambda = 1.0$: Shows near-linear scaling (better than $O(n^2)$) due to adjacency property
- $\lambda \neq 1.0$: Shows quadratic scaling as expected for exhaustive search

Table 2: Performance comparison across different lambda values on synthetic datasets

Lambda Value	N (samples)	Training Time (s) Tree / SVC	Prediction Time (s) Tree / SVC	Accuracy (%) Tree / SVC	Speedup Train / Pred
5.0	100	0.000721 / 0.001718	0.000126 / 0.004093	90.5 / 90.5	2.38x / 32.60x
	200	0.002349 / 0.003440	0.000141 / 0.013088	89.0 / 89.0	1.46x / 93.12x
	400	0.008650 / 0.009497	0.000176 / 0.045511	89.25 / 89.44	1.10x / 259.32x
	1000	0.051930 / 0.054325	0.000299 / 0.277141	89.0 / 88.92	1.05x / 927.37x
	2000	0.206405 / 0.208046	0.000427 / 1.056877	88.82 / 88.85	1.01x / 2476.05x
	5000	1.337285 / 1.428290	0.000981 / 6.594770	88.66 / 88.70	1.07x / 6719.38x
10.0	100	0.000679 / 0.001738	0.000124 / 0.004136	90.5 / 90.5	2.56x / 33.43x
	200	0.002418 / 0.003423	0.000146 / 0.013259	89.0 / 89.0	1.42x / 90.94x
	400	0.008658 / 0.009555	0.000175 / 0.045692	89.25 / 89.44	1.10x / 260.93x
	1000	0.052488 / 0.058024	0.000347 / 0.345599	89.0 / 88.92	1.11x / 996.47x
	2000	0.206459 / 0.206304	0.000449 / 1.077199	88.82 / 88.85	1.00x / 2400.72x
	5000	1.333215 / 1.413591	0.000977 / 6.619439	88.66 / 88.70	1.06x / 6778.18x
20.0	100	0.000681 / 0.001718	0.000120 / 0.004090	89.5 / 90.5	2.52x / 34.16x
	200	0.002417 / 0.003317	0.000144 / 0.013163	87.5 / 89.0	1.37x / 91.71x
	400	0.008603 / 0.009521	0.000174 / 0.045666	89.25 / 89.44	1.11x / 262.47x
	1000	0.052221 / 0.054934	0.000304 / 0.277956	89.0 / 88.92	1.05x / 913.11x
	2000	0.207620 / 0.206897	0.000438 / 1.059200	88.82 / 88.85	1.00x / 2418.84x
	5000	1.331622 / 1.415241	0.000966 / 6.633535	88.74 / 88.70	1.06x / 6867.11x
30.0	100	0.000681 / 0.001778	0.000125 / 0.004145	89.5 / 90.5	2.61x / 33.28x
	200	0.002341 / 0.003394	0.000140 / 0.013161	89.0 / 89.0	1.45x / 93.73x
	400	0.008647 / 0.009331	0.000177 / 0.045701	87.62 / 89.44	1.08x / 257.59x
	1000	0.051890 / 0.055510	0.000328 / 0.278735	89.0 / 88.92	1.07x / 851.01x
	2000	0.206779 / 0.206360	0.000439 / 1.059128	88.82 / 88.85	1.00x / 2412.61x
	5000	1.333043 / 1.411878	0.000936 / 6.603125	88.74 / 88.70	1.06x / 7054.61x

2. **Training Time Performance Comparison:** A critical observation emerges when comparing the practical performance between $\lambda = 1.0$ (Table 1) and $\lambda \neq 1.0$ (Table 2):
 - $\lambda = 1.0$ **advantage:** Achieves training speedups of 10.5x to 96.2x compared to classical SVM, with dramatic improvement as dataset size increases (N=5000: 96.2x speedup)
 - $\lambda \neq 1.0$ **limitation:** Training speedups are significantly lower, ranging from 2.4x to 1.06x, with performance actually degrading as dataset size increases (N=5000: only 1.06x speedup)
 - **Computational explanation:** The $O(n^2)$ complexity of $\lambda \neq 1.0$ cases requires exhaustive evaluation of all node pairs, negating the computational advantages of the tree structure for training. In contrast, $\lambda = 1.0$ leverages the adjacency property to achieve near-linear training complexity.
 - **Practical implication:** While $\lambda \neq 1.0$ demonstrates the theoretical generality of our method, $\lambda = 1.0$ represents the computationally optimal configuration for practical applications.
3. **Training Performance Consistency:** Despite the higher complexity for $\lambda \neq 1.0$, training times remain practical across different lambda values for each dataset size, with the tree structure still providing significant advantages over classical SVM.
4. **Prediction Performance Stability:** Prediction times show minor variations across lambda values, with all configurations maintaining excellent speedup ratios compared to classical SVM.
5. **extbfAccuracy Trade-offs:** Lower lambda values ($\lambda = 5.0, 10.0$) tend to preserve accuracy better, while higher values ($\lambda = 20.0, 30.0$) may show slight accuracy drops on some dataset sizes, particularly noticeable at N=400 and N=800.
6. **extbfOptimal Lambda Selection:** For most practical applications, $\lambda = 5.0$ or $\lambda = 10.0$ provide the best balance between computational efficiency and accuracy preservation.
7. **extbfScalability Across Parameters:** All lambda values demonstrate excellent scalability, with prediction speedups reaching over 7000x for the largest datasets regardless of the lambda choice.

6.3 Performance Analysis

The experimental results reveal several important findings:

1. **Training Speedup:** Our SVM on Tree method achieves consistent training speedups ranging from 10.5x to 96.2x compared to scikit-learn’s SVC. The speedup increases dramatically with dataset size, demonstrating exceptional scalability.
2. **Prediction Speedup:** The prediction phase shows even more dramatic improvements, with speedups ranging from 9.4x to an impressive 6595.8x for the largest dataset (N=5000). This extraordinary speedup makes our method ideal for real-time applications.
3. **Accuracy Preservation:** Our method maintains competitive accuracy across all dataset sizes, consistently matching the classical SVM performance (around 88-90% accuracy).
4. **Scalability:** The computational advantage becomes more pronounced with larger datasets, with the most significant improvements observed at N=5000 where prediction is over 6500 times faster than classical SVM.

6.4 Real-world Dataset Evaluation

We further evaluate our method on real-world datasets to demonstrate its practical applicability beyond synthetic data. For these experiments, we focus on the computationally optimal configuration with $\lambda = 1.0$, which demonstrated superior training performance in our synthetic data analysis.

6.4.1 Iris Dataset

The Iris dataset, a classic benchmark in machine learning, contains 150 samples with 4 features representing sepal and petal measurements of iris flowers. For our experiments, we applied PCA to reduce the data to 2 dimensions and converted it to a binary classification problem.

Table 3 shows the performance comparison on the Iris dataset over 10 runs:

Table 3: Performance comparison on Iris dataset (median of 10 runs)

Metric	SVM on Tree	SVC (RBF)	Speedup/Difference
Training Time	0.000052s	0.000643s	12.33x
Prediction Time	0.000101s	0.000159s	1.58x
Accuracy	100.0%	100.0%	+0.0%

The Iris dataset demonstrates excellent suitability for our tree-based approach with:

- **Suitability Score:** 5.358 (indicating high compatibility with tree structure)
- **Spine Dominance:** 211.306 (showing clear linear separability along the principal direction)
- **Zero Overlap:** 0.000 overlap on spine, indicating clean separation

Both methods achieved perfect classification accuracy (100%), but our SVM on Tree method was 12.33 times faster in training and 1.58 times faster in prediction.

6.4.2 Wine Dataset

The Wine dataset contains chemical analysis of wines from three different cultivars. We converted this to a binary classification problem and reduced the dimensionality to 2D using PCA for compatibility with our tree-based approach.

Table 4 shows the performance comparison on the Wine dataset over 10 runs:

Table 4: Performance comparison on Wine dataset (median of 10 runs)

Metric	SVM on Tree	SVC (RBF)	Speedup/Difference
Training Time	0.000075s	0.000713s	9.52x
Prediction Time	0.000101s	0.000182s	1.79x
Accuracy	88.89%	88.89%	+0.0%

The Wine dataset presents a more challenging scenario with:

- **Suitability Score:** 0.728 (moderate compatibility with tree structure)
- **Spine Dominance:** 1.885 (lower than Iris, indicating less clear linear separability)
- **Class Imbalance:** 0.331 (indicating uneven class distribution)
- **Spine Overlap:** 0.225 (some overlap along the principal direction)

Despite the more challenging characteristics of the Wine dataset, our SVM on Tree method achieved identical accuracy (88.89%) to the classical SVM while providing 9.52x speedup in training and 1.79x speedup in prediction.

6.5 Discussion

The experimental results across both synthetic and real-world datasets demonstrate the effectiveness of our SVM on Tree approach:

1. **Computational Efficiency:** Consistent speedups ranging from 1.58x to 6595.8x in prediction time and 9.52x to 96.2x in training time across different dataset sizes and types. The most impressive results are observed with larger synthetic datasets where prediction speedup exceeds 6500x.
2. **Accuracy Preservation:** Our method maintains identical accuracy to classical SVM across all tested datasets. On well-structured data like Iris, both methods achieve perfect classification (100%), while on more complex datasets like Wine, both achieve competitive performance (88.89%).
3. **Scalability:** The performance advantage becomes more pronounced with larger datasets, making our approach particularly suitable for big data applications. Synthetic experiments show speedups exceeding 6500x for prediction on the largest datasets (N=5000).
4. **Dataset Adaptability:** The method works effectively across datasets with varying characteristics:
 - High suitability (Iris: score 5.358) with perfect linear separability
 - Moderate suitability (Wine: score 0.728) with class imbalance and spine overlap
 - Synthetic datasets with controlled noise and separation parameters
5. **Practical Applicability:** The method works effectively on real-world datasets, not just synthetic ones, demonstrating its practical value for machine learning applications.

These results validate our theoretical analysis and confirm that the tree-based formulation provides a computationally efficient alternative to classical SVM without sacrificing classification quality. The suitability score metric effectively predicts when our method will perform optimally, helping practitioners decide when to apply this approach.

The synthetic results demonstrate that our SVM on Tree method successfully addresses the computational challenges of traditional SVM while preserving classification quality. The significant speedups in both training and prediction phases, combined with maintained accuracy, make our approach particularly suitable for large-scale applications and real-time systems.