# Fast SVM on Tree

Cong-Huan Tran[1], Thu-Le Tran[*] [2],
Khanh-Duy Le[2], Ngan Nguyen[2], and Clément Elvira[3]

[1]Faculty of xxx, College of Information and Communication Technology, Cantho
University, Can Tho City, Vietnam
[2]Faculty of Mathematics and Informatics, Teacher College,, Cantho University, Can
Tho City, Vietnam
[3]France

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

---

[*]Corresponding author, ttle@ctu.edu.vn

# 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 $i = 1, ..., n$

2. Classical soft-margin SVM model

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

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 $i = 1, ..., n$

2. Denote $V_+ = \{x_i \in V : y_i = +1\}, V_- = \{x_i \in V : y_i = -1\}$

3. Let $\Delta$ be the line passing through the centers of $V_+$ and $V_-$

4. Let $x_i'$ be the projection of $x_i$ onto $\Delta$

5. Tree $\mathcal{T} = (V, E)$ be a tree with vertex set $V = \{x_1, ...., x_m, x_1', ..., x_n'\} \subset \mathbb{R}^n$ and edge set $E = \{(x_i', x_{i+1}') : i = 1, ..., n - 1\} \cup \{(x_i, x_i') : i = 1, ..., n\}$

6. The SVM on tree aims to minimize the objective function

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

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 writting.]

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

$$w^T x + b \geq 1, \quad \forall y_i = +1$$
$$w^T x + b \leq -1, \quad \forall y_i = -1$$

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:

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

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.,

$$w^s + b = 1$$
$$w^p + b = -1.$$

Hence the name SVM.

## 3.2 Soft Margin SVM

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

$$\min_{w,b,\xi} \quad \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} \xi_i \tag{3}$$
$$\text{s.t.} \quad y_i(w^T x_i + b) \geq 1 - \xi_i, \quad i = 1, \ldots, n \tag{4}$$
$$\xi_i \geq 0, \quad i = 1, \ldots, n \tag{5}$$

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, Section 5.5 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 $\{(x_i, y_i) \in \mathbb{R}^d \times \{-1, +1\} : i = 1, \ldots, n\}$, we describe a procedure for constructing a tree that approximates the dataset.

Let $\Delta$ be a line in $\mathbb{R}^d$. In this work, we focus on the line $\Delta$ connecting the two class means. For each $i = 1, \ldots, n$, let $x_i'$ denote the orthogonal projection of $x_i$ onto $\Delta$. 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 $\Delta$, 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:
$$\mu_- = \frac{1}{|V_-|} \sum_{y_i = -1} x_i,$$
$$\mu_+ = \frac{1}{|V_+|} \sum_{y_i = +1} x_i.$$

The unit vector pointing the direction of the line $\Delta$ conencting 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 $\mu_-$ in direction $w$

$$x_i' = \mu_- + t_i w,$$
$$t_i = \langle x_i - \mu_-, \, w \rangle.$$

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

$$t_1 \leq \cdots \leq t_n.$$

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

$$V = \{x_1, \ldots, x_n\} \cup \{x_1', \ldots, x_n'\},$$

and edge set

$$E = \{(x_i', x_{i+1}') : i = 1, \ldots, n-1\} \cup \{(x_i, x_i') : i = 1, \ldots, n\}.$$

**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 <u>spine edges</u>, and to edges connecting each original point with its projection as <u>spoke edges</u>. The path formed by the sequence of all spine edges is called the <u>spine</u> of the tree.
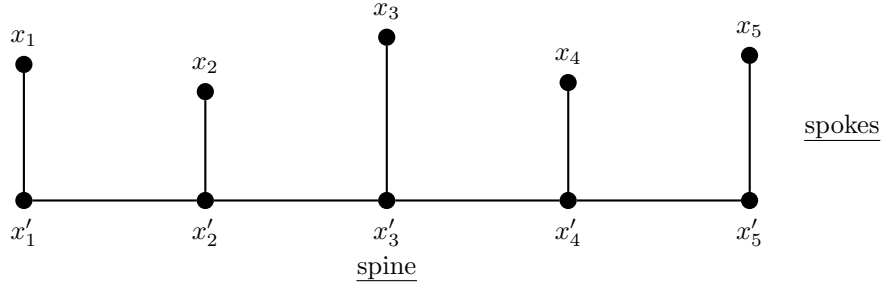


Figure 1: Augmented tree structure showing original data points $\{x_i\}$ and their projections $\{x_i'\}$ onto the spine. The spine connects consecutive projections, while spoke edges connect each original point to its projection.

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 $i = 1, \ldots, n$. 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:

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

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

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]

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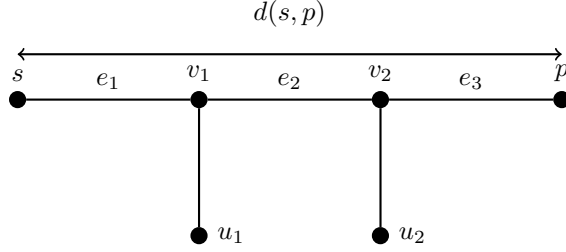
Figure 2: Support pair illustration showing nodes $s$ and $p$ connected by a path with intermediate vertices and subtree nodes.

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 $O(n^3)$.

However, we can significantly improve this complexity by leveraging the tree structure and adjacency properties.

**Theorem 5.1.** *For any $\lambda \geq 0$, we can find the support nodes in $O(n^2)$.*

**Proof idea:** The key insight is that while the distance matrix on tree can be computed in $O(n^2)$ time, the noise function $f(u, v)$ can also be computed in $O(n^2)$ time for all pairs by leveraging dynamic programming on the tree structure. This eliminates the need for $O(n^3)$ pairwise evaluations.

While the general case admits an $O(n^2)$ algorithm, we can achieve even more significant improvements by focusing on the special case $\lambda = 1$. This parameter choice leads to a remarkable structural property that dramatically reduces the search space and enables linear-time preprocessing algorithms.

## 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 $x_1', x_2', \ldots, x_n'$ 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.,

$$d(s, p_1) = d(s, p_2) + d(p_2, p_1) \tag{8}$$

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

8

extend the boundary from $p_2$ to $p_1$. This set is crucial for analyzing the difference in noise functions between different support pairs. [LE: Added formal definition of residual set concept]

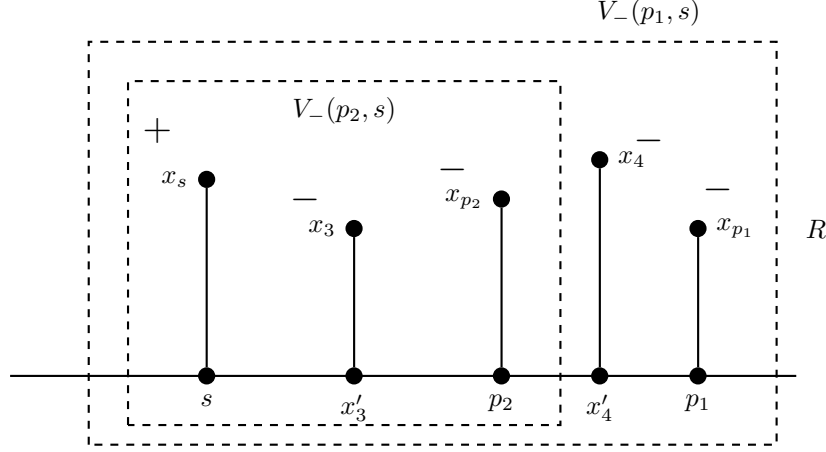Please refer to Figure 3 for an illustration of the difference set $R$.



Figure 3: Difference set $R$ illustration.

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

**Lemma 5.2.** *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

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

9

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 5.3** (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. □



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

The adjacency property established in Theorem 5.3 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.

## 5.5  Algorithm and Complexity Analysis

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

### 5.5.1 Preprocessing

Recall $x'_1, \ldots, x'_n$ are projections of $x_1, \ldots, x_n$ and we assume that $x'_1, \ldots, x'_n$ have been sorted in a certain order, e.g., by their first coordinates. By the adjacency property (Theorem 5.3), we know that support vectors form a pair $(x'_i, x'_{i+1})$ with opposite signs. To compute the total loss, we only need $f(x'_i, x'_{i+1})$ for $i = 1, .., n-1$. In the following, we will clarify how to compute these values efficiently using dynamic programming.

For $i = 1, \ldots, n-1$, we define data at iteration $i$ as

$$
\begin{aligned}
N_+^{>i} &= |\{v \in V_+(x'_i, x'_{i+1})\}| = 2|\{j > i : y_j = +1\}| \\
N_-^{>i} &= |\{v \in V_-(x'_i, x'_{i+1})\}| = 2|\{j > i : y_j = -1\}| \\
N_+^{\leq i} &= |\{v \in V_+(x'_{i+1}, x'_i)\}| = 2|\{j \leq i : y_j = +1\}| \\
N_-^{\leq i} &= |\{v \in V_+(x'_{i+1}, x'_i)\}| = 2|\{j \leq i : y_j = -1\}|
\end{aligned}
$$

and

$$
\begin{aligned}
d_+^{>i} &= \sum_{v \in V_+(x'_i, x'_{i+1})} d(x'_i, v) & &= \sum_{j > i, y_j = +1} [d(x'_i, x'_j) + d(x'_i, x_j)], \\
d_-^{>i} &= \sum_{v \in V_-(x'_i, x'_{i+1})} d(x'_i, v) & &= \sum_{j > i, y_j = -1} [d(x'_i, x'_j) + d(x'_i, x_j)], \\
d_+^{\leq i} &= \sum_{v \in V_+(x'_{i+1}, x'_i)} d(v, x'_{i+1}) & &= \sum_{j \leq i, y_j = +1} [d(x'_j, x'_{i+1}) + d(x_j, x'_{i+1})] \\
d_-^{\leq i} &= \sum_{v \in V_-(x'_{i+1}, x'_i)} d(v, x'_{i+1}) & &= \sum_{j \leq i, y_j = -1} [d(x'_j, x'_{i+1}) + d(x_j, x'_{i+1})].
\end{aligned}
$$

Consider $i = 1, \ldots, n-1$. If $y_i = -1$ and $y_{i+1} = +1$, we then have

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

Otherwise, $y_i = +1$ and $y_{i+1} = -1$, then

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

Therefore, to compute the total noise function, we should have an efficient method to compute $(d_+^{>i}, d_-^{>i}, d_+^{\leq i}, d_-^{\leq i})$. This can be achived using the dynamic updates.

The following proposition claims that how the information of $N_-^{\leq i}$ and $d_-^{\leq i}$ can be used to derive $N_-^{\leq i+1}$ and $d_-^{\leq i+1}$ efficiently. The proof of the proposition will show the procedure in details.

**Proposition 5.4.** *For $i = 1, \ldots, n-1$, if we have $N_-^{\leq i}$ and $d_-^{\leq i}$, then we can update $N_-^{\leq i+1}$ and $d_-^{\leq i+1}$ in constant time.*

*Proof.* By definition,

$$
d_-^{\leq i} = \sum_{\substack{j \leq i \\ y_j = -1}} \big(d(x'_j, x'_{i+1}) + d(x_j, x'_{i+1})\big), \qquad d_-^{\leq i+1} = \sum_{\substack{j \leq i+1 \\ y_j = -1}} \big(d(x'_j, x'_{i+2}) + d(x_j, x'_{i+2})\big).
$$

First consider the sum over indices $j \leq i$ with $y_j = -1$. Due to the tree structure consisting of the spine (connecting the $x'_k$ in order) and spokes (connecting each $x_k$ to its projection $x'_k$), for all $j \leq i$, the unique path from $x'_j$ (or $x_j$) to $x'_{i+2}$ passes through $x'_{i+1}$. Therefore,

$$
d(x'_j, x'_{i+2}) = d(x'_j, x'_{i+1}) + d(x'_{i+1}, x'_{i+2}), \qquad d(x_j, x'_{i+2}) = d(x_j, x'_{i+1}) + d(x'_{i+1}, x'_{i+2}).
$$

Summing over all $j \leq i$ with $y_j = -1$, the increase due to shifting the anchor from $x'_{i+1}$ to $x'_{i+2}$ is

$$N_-^{\leq i} \cdot d(x'_{i+1}, x'_{i+2}),$$

since $N_-^{\leq i}$ counts exactly the number of negative vertices participating in the sum (each index $j$ contributes two vertices $x'_j$ and $x_j$).

Additionally, if $y_{i+1} = -1$, two new vertices appear in the region $\leq i+1$, namely $x'_{i+1}$ and $x_{i+1}$. Their contribution to $d_-^{\leq i+1}$ is

$$d(x'_{i+1}, x'_{i+2}) + d(x_{i+1}, x'_{i+2}) = d(x'_{i+1}, x'_{i+2}) + \big(d(x_{i+1}, x'_{i+1}) + d(x'_{i+1}, x'_{i+2})\big) = 2\, d(x'_{i+1}, x'_{i+2}) + d(x_{i+1}, x'_{i+1}).$$

Combining both parts, we obtain the update formulas for the two cases based on the label $y_{i+1}$:

$$N_-^{\leq i+1} = \begin{cases} N_-^{\leq i} + 2, & \text{if } y_{i+1} = -1, \\ N_-^{\leq i}, & \text{if } y_{i+1} = +1, \end{cases}$$

$$d_-^{\leq i+1} = \begin{cases} d_-^{\leq i} + N_-^{\leq i}\, d(x'_{i+1}, x'_{i+2}) + 2\, d(x'_{i+1}, x'_{i+2}) + d(x_{i+1}, x'_{i+1}), & \text{if } y_{i+1} = -1, \\ d_-^{\leq i} + N_-^{\leq i}\, d(x'_{i+1}, x'_{i+2}), & \text{if } y_{i+1} = +1. \end{cases}$$

Clearly, each update step only uses quantities known at step $i$ and distances between consecutive spine points, so it can be performed in constant time. Therefore, knowing $N_-^{\leq i}$ and $d_-^{\leq i}$ is sufficient to update $N_-^{\leq i+1}$ and $d_-^{\leq i+1}$. $\qquad\square$

Extending above proposition, we obtain the following result.

**Proposition 5.5.** *We can compute the following quantities in $O(n)$.*

1. *$(N_+^{>i}, N_-^{>i}, N_+^{\leq i}, N_-^{\leq i})$ for all $i = 1, \ldots, n-1$*

2. *$(d_+^{>i}, d_-^{>i}, d_+^{\leq i}, d_-^{\leq i})$ for all $i = 1, \ldots, n-1$*

3. *$f(x'_i, x'_{i+1})$ for all $i = 1, \ldots, n-1$.*

*Proof.* We prove this by applying Proposition 5.1 systematically, using the dynamic programming technique on trees.

*Step 1: Computing all counting quantities in $O(n)$.* The counting quantities can be computed by simple linear scans:

- $N_+^{>i}, N_-^{>i}$: Scan from right to left, counting positive/negative labels.

- $N_+^{\leq i}, N_-^{\leq i}$: Scan from left to right, counting positive/negative labels.

Each quantity is updated in $O(1)$ per position, giving total $O(n)$ time.

*Step 2: Computing all distance quantities in $O(n)$.* This is where we apply the dynamic programming approach described above. The key insight is that our augmented tree has a spine structure, which allows efficient re-rooting operations.

Left-to-right pass for computing $d_+^{\leq i}, d_-^{\leq i}$:

- Initialize: $d_+^{\leq 1} = 2d(x'_1, x'_2) + d(x_1, x'_2)$ if $y_1 = +1$, else $0$.

- For $i = 1, \ldots, n-2$: Apply Proposition 1 (and its symmetric version for positive labels) to update $d_\pm^{\leq i+1}$ from $d_\pm^{\leq i}$ in $O(1)$ time.

Right-to-left pass for computing $d_+^{>i}, d_-^{>i}$:

- Initialize: $d_+^{>n-1} = 0, d_-^{>n-1} = 0$ (no points beyond position $n-1$).

- For $i = n-1, n-2, \ldots, 2$: Use the symmetric update formula to compute $d_\pm^{>(i-1)}$ from $d_\pm^{\geq i}$ in $O(1)$ time.

The symmetric update formula for the right-to-left pass is obtained by reversing the roles: when moving from anchor $x'_{i+1}$ to $x'_i$, all points $j > i$ have their distances increased by $d(x'_i, x'_{i+1})$, and if $y_i = -1$, we add the contribution of the new point $i$.

*Step 3: Computing $f(x'_i, x'_{i+1})$ in $O(n)$.* Once all counting and distance quantities are computed, we can evaluate $f(x'_i, x'_{i+1})$ for each $i = 1, \ldots, n-1$ using the formulas given in the preprocessing section. Each evaluation takes $O(1)$ time, giving total $O(n)$ time.

$\square$

### 5.5.2 Algorithm and Complexity Analysis

---
**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 $\mu_-$, $\mu_+$ and spine direction $w$
4: Project all points: $x'_i = \mu_- + \langle x_i - \mu_-, w \rangle w$ for $i = 1, \ldots, n$
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 $i = 1, \ldots, n-1$ using preprocessing step in Section 5.5.1
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 $i = 1, \ldots, n-1$
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** $(s^*, p^*)$

---

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.

- **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 5.3.

We finally obtain

**Theorem 5.6** (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 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.

Table 1: Performance comparison on synthetic datasets

| N (samples) | Training Time (s) Tree / SVC | Prediction Time (s) Tree / SVC | Accuracy (%) Tree / SVC | Speedup 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 |

## 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:

- Spine separation: $sep = 6.0$

- 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.

## 6.2 Synthetic Dataset Results

Table 1 shows the performance comparison on synthetic datasets of varying sizes. The results demonstrate consistent speedups in both training (FIT) and prediction (PRED) phases while maintaining competitive accuracy.

## 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.

### 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 2 shows the performance comparison on the Iris dataset over 10 runs:

Table 2: 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 3 shows the performance comparison on the Wine dataset over 10 runs:

Table 3: 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.