

Application of Geometric Graph Distances in Machine Learning Classification

William Rodman

April 2024

Abstract

Geometric graph theory, prevalent in computer science for structuring networks like the World Wide Web and digital road maps, lacks computationally efficient methods for measuring distances between graphs. This thesis investigates the traversal distance algorithm's potential in machine learning classification, specifically in enhancing the precision of the k-nearest neighbors (k-NN) model for classifying geometric graphs represented as English letters. Prior research has explored various algorithms for this purpose, the hypothesis being, the traversal distance offers computational efficiency despite its asymmetry. Employing a dataset of English letters, the research tests the hypothesis that a k-NN model, integrated with traversal distance, could achieve high classification precision. Through comparative analysis with the graph edit distance (GED), the traversal distance's effectiveness in classification tasks is established. Demonstrating the traversal distance's potential application in supervised machine learning.

Contents

Acknowledgments	4
1 Introduction	5
2 Introducing the Traversal Distance	6
2.1 Properties of Geometric Graph Theory	6
2.2 Weak Fréchet Distance	7
2.3 Free-Space	9
2.4 Traversal Distance and Algorithm	10
2.5 Properties of the Traversal Distance	13
2.6 Geometric Graph Edit Distance	14
3 Visualizing the Traversal Distance	16
3.1 Weak Fréchet Distance Free-Space Diagram	16
3.2 Traversal Distance Free-Space Visualization	18
3.3 Example of Traversal Distance Visualization	21
4 Distance Measurements in Machine Learning	24
4.1 Euclidean Distance Between Two Points	24
4.2 Introducing the K-Nearest Neighbors Model	24
4.3 Evaluating K-Nearest Neighbors Predictions	26
5 Applying the Traversal Distance to Classification Problems	28
5.1 Symmetric Case of the Traversal Distance	28
5.2 K-Nearest Neighbors Using the Traversal Distance	29
5.3 K-Nearest Neighbors Using Graph Edit Distance	31
6 Conclusion	32
Appendix	33
References	36

Acknowledgments

I want to extend my gratitude to the individuals who have guided and supported me throughout my research. Their expertise, encouragement, and mentorship played a important role in shaping this thesis.

First and foremost, thank you to Dr. Carola Wenk, Professor at Tulane University's Department of Computer Science. Dr. Wenk introduced me to research in my freshman year of college as a research assistant. Under her National Science Foundation research grant, "A Unified Framework for Geometric and Topological Signature-Based Shape Comparison," I have been able to gain experience in research over the last three years. This research laid the foundation for my honors thesis, where I have continued to work with Dr. Wenk as my primary thesis advisor. Next, I would like to thank Erfan Hosseini, a Computer Science PhD candidate at Tulane University, advised by Dr. Wenk. Erfan has helped me countless times, both in my time as a research assistant and throughout the research for my honors thesis. Thank you to Dr. Ramgopal Mettu, Professor in the Department of Computer Science at Tulane University, for his role as my second thesis reader. I also extend my thanks to Dr. Sushovan Majhi, a Post-doc researcher at the University of California, Berkeley. Working with Dr. Majhi has deepened my understanding of the Graph Edit Distance. His contribution of the English letter dataset was an important component of this thesis. My gratitude also goes to Dr. Liz Munch and Dr. Sarah Percival, faculty at Michigan State University, for their advice and provision of the plant leaf dataset to the research group at Tulane University.

In addition to my academic mentors, the support of my family has also been a source of motivation. A special thank you to my father, Rick Rodman, who supported me throughout my bachelor's degree. His career at NASA Goddard Space Flight Center sparked my interest in mathematics and computer science. To my brother, Sander Rodman, an undergraduate at Cornell University studying information and data science. And finally, to my mother and stepfather, Laura and Todd Ecker.

1 Introduction

Geometric graph theory has become a significant concept in modern mathematics. Its relevance extends into the field of computer science, where it plays a crucial role in defining and structuring complex networks, including the World Wide Web and digital road mapping systems like Apple Maps [2]. The widespread application of geometric graphs in computational geometry has garnered interest from research organizations, including the National Science Foundation. This interest has led to funding for academic research focused on advancing methods for comparing, measuring, and efficiently storing geometric graphs.

The measurement of distance between two geometric graphs is a practical comparison in geometric graph theory. However, the comparison faces challenges since there exists no closed-form solutions for computing the spatial distance between graphs. This challenge has prompted the development of a variety of algorithms designed to approximate this distance. The focus of this thesis is to investigate the traversal distance algorithm, examining its advantages, disadvantages, and applications [1]. One advantage of this algorithm is its computational efficiency; it has a polynomial time complexity, contrasting the NP-Hard complexity of the graph edit distance (GED). However, a notable limitation is its asymmetry, an issue that GED avoids [9].

This thesis will test the application of the traversal distance in the classification of geometric graphs, a challenge within supervised machine learning. Supervised learning, a subfield of computer science and artificial intelligence, has gained significant traction, particularly following OpenAI's release of ChatGPT [6]. The central hypothesis of this research is that a k-nearest neighbors (k-NN) machine learning model, when integrated with the traversal distance as its distance metric, will prove precise in classifying geometric graphs. This hypothesis will be empirically tested using a dataset of English letters represented as geometric graphs.

2 Introducing the Traversal Distance

In this chapter, the discussion begins by explaining what geometric graphs are; the inputs needed to compute the traversal distance. Next, the chapter introduces a simpler way to measure the distance between two curves, known as the weak Fréchet distance. This topic acts as a stepping stone, helping to better understand the more complex traversal distance. Finally, the advantages and disadvantages of using the traversal distance are discussed. It is compared with other methods of measuring distances in geometric graphs, like the graph edit distance (GED), to gauge its usefulness.

2.1 Properties of Geometric Graph Theory

Geometric graph theory exists within the broader field of graph theory, focused on the study of graphs embedded in a geometric space. In the case of this thesis, graphs are embedded in Euclidean planes, where vertices represent coordinate points and edges signify line segments between these points.

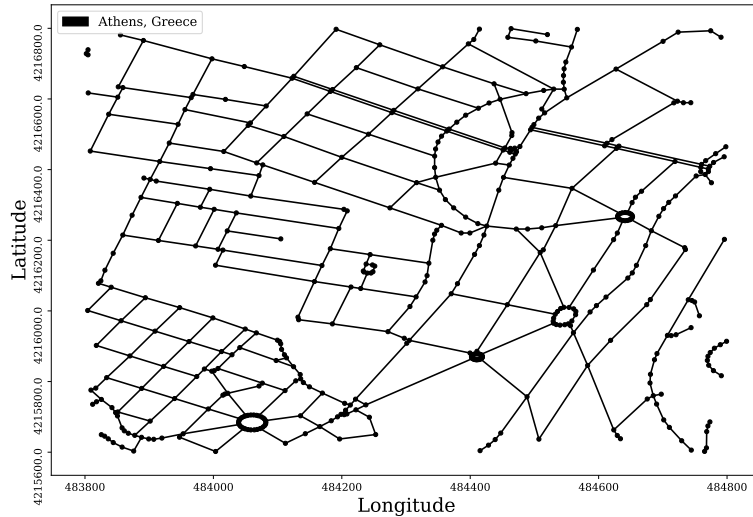


Figure 1: Geometric graph capturing network of roads in Athens, Greece [7].

This discipline extends into numerous applications in computer science. Consider a digital road map, a familiar tool used while driving a car. In this context, the roads can be thought of as the edges of a geometric graph, while the intersections where these roads meet are the vertices. While traditional road maps reference geographic coordinates to pinpoint locations, digital road maps, when

stored as geometric graphs, contain embedded coordinates that correspond to locations.

Definition 1. *Let a geometric graph G be defined as a pair $G = (V, E)$, where V is the set of vertices, with each vertex $v \in V$ corresponding to a pair of two-dimensional coordinates (x, y) in the Euclidean plane. Then the set V is defined as:*

$$V = \{1, 2, \dots, n\} \quad \text{where } n \in \mathbb{N} \quad \text{such that } v_i \mapsto (x_i, y_i)$$

E is the set of edges, where each edge $e \in E$ defines a line segment bound by two vertices. An edge e comprised of vertices v_i and v_j is denoted as e_i . Then the set E is defined as:

$$E = \{e_i | (v_i, v_j) \in V \times V\} \quad \text{such that } e_i = (v_i, v_j)$$

In this structure, the geometric graph G exists entirely in the Euclidean plane [11]. The implementation of this geometric graph within a Python program takes on a slightly different form. The Python program in this thesis stores geometric graphs using two dictionaries. The first dictionary, which will be referred to as **nodes**, contains all the vertices of the graph. In this **nodes** dictionary, the keys are used for vertex identification, and the values are tuples containing the coordinates of each vertex. The **nodes** dictionary can be written as:

$$\text{nodes} = \{1 : (x_1, y_1), 2 : (x_2, y_2), \dots, n : (x_n, y_n)\}$$

The second dictionary, **nodeLinks**, stores all information pertaining to the edges of the geometric graph. In this dictionary, the keys contain vertices that form part of an edge's line segment. The values associated with these keys are lists, each containing the neighboring vertices that each completes an edge's line segment. The **nodeLinks** dictionary can be written as [15]:

$$\text{nodeLinks} = \{1 : \{j \in e_1\}, 2 : \{j \in e_2\}, \dots, n : \{j \in e_n\}\}$$

The configuration of geometric graphs in the Python program, as detailed in a later section of this chapter, reduces the time for searching for all nodes connected to a given node. This improvement in search time, however, increases the space required for storing geometric graphs [11].

2.2 Weak Fréchet Distance

Before mentioning the traversal distance, it is instructive to first understand the weak Fréchet distance, a metric for measuring the distance between two curves. The rationale for starting with the weak Fréchet distance is that it presents a simpler problem, which provides a foundation for grasping the more generalized traversal distance problem.

Letting C_1 and C_2 be defined as two polygonal curves, where each curve has a first and last point. Essentially, the weak Fréchet distance is a value ϵ , where

ϵ is the parameter for free-space spanning the first and last points of C_1 and C_2 . More strictly, the value ϵ is a weak Fréchet distance if and only if there exists a continuous path across the free-space, such that the path starts at the free-space for the first points of C_1 and C_2 and ends at the free-space for the last points of C_1 and C_2 .

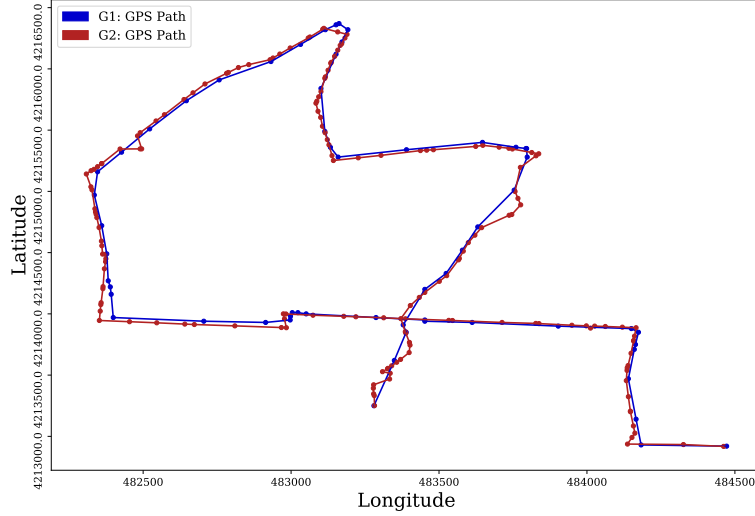


Figure 2: GPS coordinate paths as polygonal curves [16].

To illustrate, consider the scenario in Figure 2 comparing a GPS-tracked route to a hiking trail stored digitally. Here, the weak Fréchet distance would be used to determine the greatest deviation of the GPS route from the hiking trail, with the GPS trajectory and the hiking trail representing the two curves in question.

Definition 2. *Given two curves C_1 and C_2 that exist in the Euclidean plane. Let $C_1 : I = [0, 1] \rightarrow \mathbb{R}^2$, $C_2 : J = [0, 1] \rightarrow \mathbb{R}^2$, and $\|\cdot\|$ denote the Euclidean norm. Then the weak Fréchet distance $\delta_F(C_1, C_2)$ is defined as:*

$$\delta_F(C_1, C_2) := \inf_{\alpha \rightarrow [0,1], \beta \rightarrow [0,1]} \max_{t \in [0,1]} \|C_1(\alpha(t)) - C_2(\beta(t))\|,$$

where $\alpha : [0, 1] \rightarrow I$ and $\beta : [0, 1] \rightarrow J$ range over continuous parametrizations with $\alpha(0) = \text{start}_I$, $\alpha(1) = \text{end}_I$, $\beta(0) = \text{start}_J$, and $\beta(1) = \text{end}_J$.

This introduction of the weak Frechet distance equation sets the stage for the definition of the traversal distance in section 2.4 [1].

2.3 Free-Space

Before to diving into the traversal distance and algorithm, it is necessary to introduce the data structure used for storing the parameter space between two polygonal curves. In this thesis, this data structure will be referred to as the free-space FS_ϵ . FS_ϵ represents the parameter space of all possible edge and vertex combinations for two polygonal curves, where ϵ is an arbitrary polygonal threshold within the parameter space separating C_1 and C_2 so long as $\epsilon > 0$. Each combination of edges within C_1 and C_2 is designated a square free-space cell. This concept aligns with Definition 2, where α and β span the continuous parameter space of C_1 and C_2 [1].

Figure 3 shows an example of a pair of edges, e_i and e_j , taken from the polygonal curves C_1 and C_2 , respectively. Figure 4 illustrates the free-space cell of e_i and e_j , highlighting the parameter space in white when $\epsilon = 6$ [15].

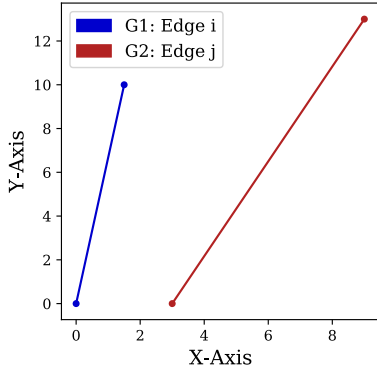


Figure 3: Pair of distinct edges [15].

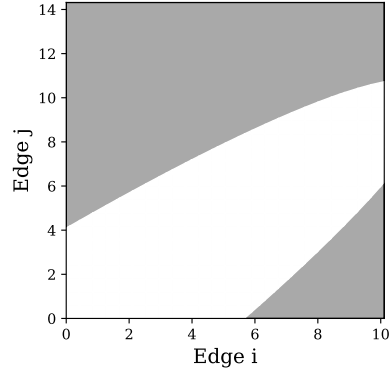


Figure 4: Corresponding free-space cell when $\epsilon = 6$ [15].

While Definition 2 describes the parametrization of space with curves, the weak Fréchet distance algorithm actually represents FS_ϵ using discrete space [1]. In this case, a free-space cell consists of four cell boundaries, each representing a wall of the square free-space cell. Each cell boundary within a free-space cell stores the starting and ending points, denoted $(start, end)$, of free-space along a cell wall. These eight points collectively form a polygon that contains the cell's free-space [15].

Using the same pair of edges from Figure 3, Figure 5 illustrates a cell boundary calculation for an edge and vertex combination. Figure 6 illustrates how the cell boundary calculation from Figure 5 is stored inside the corresponding free-space cell

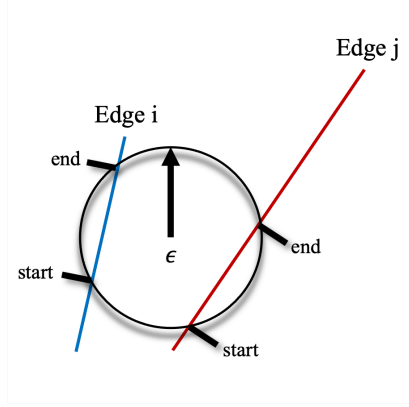


Figure 5: Cell boundary starting and ending points calculated based on value ϵ .

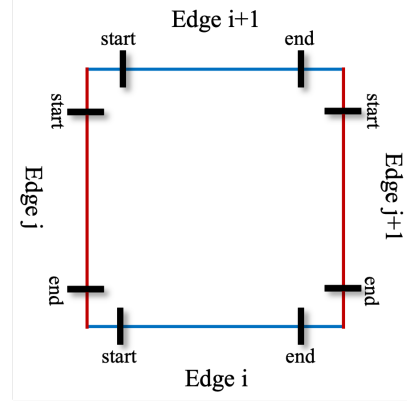


Figure 6: Cell boundary starting and ending points inside a free-space cell.

As will be shown in the next section, this explanation of free-space and the use of cell boundaries applies to two geometric graphs G_1 and G_2 [1]. The Python program stores `cell_boundaries` in the form of a dictionary [15]:

$$cell_boundaries = \{(e, v) : cellBoundary(e, v) \mid e \in E_1, E_2 \quad v \in V_1, V_2\}$$

$$cellBoundary(e_i, v_j) = (start, end) \quad \text{where} \quad 0 \leq start < end \leq 1$$

2.4 Traversal Distance and Algorithm

The traversal distance should be thought of as a more general form of the weak Fréchet distance. While the weak Fréchet distance is concerned with measuring the distance between two curves C_1 and C_2 , the traversal distance extends this concept to measure the distance between two geometric graphs G_1 and G_2 .

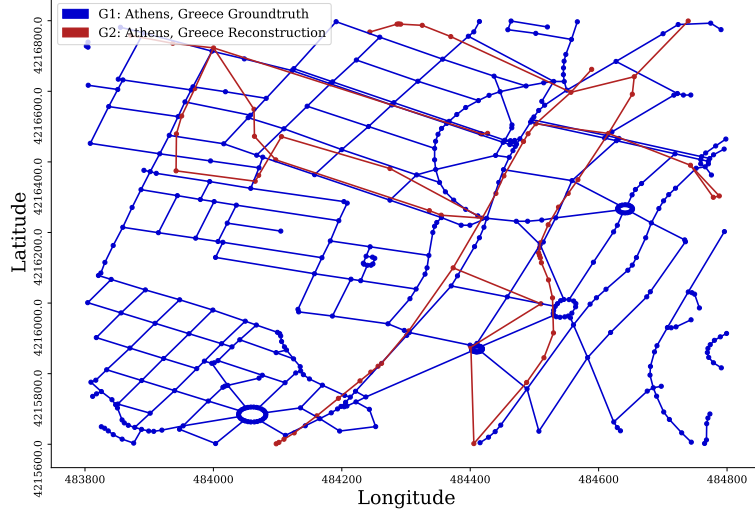


Figure 7: Network of roads in Athens, Greece and constructed GPS paths as geometric graphs [7].

To understand this, consider the traversal distance from G_1 to G_2 as the maximum distance required to cover any point on the edges of G_1 , while simultaneously traversing the entirety of G_2 in a continuous fashion. A key distinction between the traversal distance and the weak Fréchet distance is their symmetry properties. The weak Fréchet distance is symmetric, meaning the distance from C_1 to C_2 is identical to that from C_2 to C_1 , as it completely traverses both curves. In contrast, the traversal distance is asymmetric, such that the distance from G_1 to G_2 may differ from the distance from G_2 to G_1 , since the entire traversal of G_2 in a continuous fashion is not required.

Definition 3. Given two geometric graphs $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$. Define the traversal distance δ_T from G_2 to G_1 as:

$$\delta_T(G_1, G_2) = \inf_{f, g} \max_{t \in [0, 1]} \|f(t) - g(t)\|$$

where f traverses entirely over G_2 and g traverses partially over G_1 .

After defining the traversal distance equation, it is important to note that due to its nature as an infimum, the traversal distance algorithm is computationally expensive. This characteristic creates an opportunity for a less expensive algorithm for approximating the traversal distance, offering a practical solution for evaluating this metric in real-world applications [1].

The algorithm designed to decide traversal distance is comprised of three components. The first component, a dynamic program, populates `cell_boundaries`

for a given ϵ . The second component verifies that G_2 has been completely traversed by the first algorithm. This is done by projecting `cell_boundaries` onto G_2 , if the projection covers the entirety of G_2 , then the verification is true. When verification is true, these two algorithms are responsible for deciding whether $\delta_T(G_1, G_2) \leq \epsilon$. To approximate the infimum of the traversal distance, a binary search algorithm is implemented. This algorithm minimizes the distance to meet specific precision criteria, ultimately yielding an approximation ϵ where $\delta_T(G_1, G_2) \leq \epsilon \wedge \delta_T(G_1, G_2) \approx \epsilon$ [15].

Step 1: Compute the Cell Boundaries. Defined in Algorithm 1, the dynamic program used to compute `cell_boundaries` is a Depth-First Search (DFS) algorithm `DFSTraversalDist`. This algorithm writes to `cell_boundaries` while traversing all cell boundaries in FS_ϵ .

Algorithm 1 `DFSTraversalDist`

```

Initialize an empty set visited
Initialize an empty dictionary cell_boundaries

procedure DFSTRAVERSALDIST(CB)
  Compute CB
  Add CB to visited
  Insert CB in cell_boundaries

  for each neighbor of CB in nodeLink do
    if neighbor is not in visited then
      DFSTRAVERSALDIST(CB)
    end if
  end for
end procedure

seed_CB = cellBoundary( $e_1, v_1$ )                                 $\triangleright$  Starting cell boundary.
Call DFSTRAVERSALDIST(seed_CB)

```

Step 2: Verify the Traversal of G_2 . Following the traversal detailed in **Step 1**, it is necessary to next verify that G_2 was entirely traversed. As previously discussed, this verification `projection_check` is determined true if the projection of `cell_boundaries` onto G_2 covers the entirety of the graph and is determined false otherwise. Furthermore, it is implied that $\delta_T(G_1, G_2) \leq \epsilon$ when `projection_check` equals True and $\delta_T(G_1, G_2) > \epsilon$ when `projection_check` equals False. Let the output of `projection_check` be:

$$projection_check(cell_boundaries, G_2) = \begin{cases} \text{True} & \text{Projection covers } G_2 \text{ entirely.} \\ \text{False} & \text{Otherwise.} \end{cases}$$

Step 3: Binary Search for Traversal Distance. The algorithms established in **Step 1** and **Step 2** will now be incorporated into the binary search algorithm, denoted as **binarySearch**. Defined in Algorithm 2, this particular algorithm approximates the infimum of the traversal distance equation. It achieves this by searching for the smallest value of ϵ , for which **projection_check** yields true.

Algorithm 2 **binarySearch**

```

procedure BINARYSEARCH(left, right, precision)
    Initialize  $\epsilon$ 

    while right − left > precision do
         $\epsilon \leftarrow (left + right)/2$ 
        Initialize cell_boundaries
        Call DFSTraversalDist(CB)

        if projection_check(cell_boundaries,  $G_2$ ) is True then
            right  $\leftarrow \epsilon$ 
        else
            left  $\leftarrow \epsilon$ 

        end if
    end while
    return right
end procedure

```

Given the continuous domain of the search space for ϵ , since $\delta_T(G_1, G_2) \in [0, \infty)$, it is necessary to bound ϵ within a finite search domain. Thus, we assume $\delta_T(G_1, G_2) \in [left, right]$, ensuring that the **binarySearch** is contained within $[\epsilon - precision, \epsilon]$. Here, **left** and **right**, respectively, denote the lower and upper boundaries of the search space, while **precision** specifies the degree of accuracy to which the value of ϵ is returned.

Having explained the steps of the traversal distance algorithm, it is now evident that the traversal distance is computed in the Python program by calling the **binarySearch** function.

2.5 Properties of the Traversal Distance

Space Complexity of the Cell Boundaries After running the traversal distance algorithm, the program stores the values of ϵ and **cell_boundaries**. To determine the program memory requirements for this, we calculate the upper bound of the number of cell boundaries that could exist between two geometric graphs. This calculation reveals that the space complexity of **cell_boundaries**, given two geometric graphs G_1 and G_2 , is asymptotically bound by [15]:

$$S \in O((|V_1| \times |E_2|) + (|V_2| \times |E_1|))$$

Time Complexity of the Traversal Distance The time complexity of the traversal distance algorithm is determined by the combined time complexities of the `DFSTraversalDist`, `projection_check`, and `binarySearch` algorithms.

First, consider the fact that `DFSTraversalDist` is, by definition, a DFS algorithm, which runs in polynomial time. If we assume the time it takes to compute `cellBoundry` $\in O(1)$ then the time complexity of `DFSTraversalDist` is:

$$T_D \in O((|V_1| + |E_1|) \times (|V_2| + |E_2|))$$

The `projection_check` runs in polynomial time such that [1]:

$$T_P \in O((|V_1| \times |E_2|) + (|V_2| \times |E_1|))$$

For `binarySearch` operating over a continuous space, its time complexity is influenced by the number of iterations required to achieve the desired precision within the defined bounds. Given the `left` bound, `right` bound, and `precision`, the time complexity can be articulated as [8]:

$$T_B \in O(\log_2 \Delta) \quad \text{where} \quad \Delta = \frac{\text{right} - \text{left}}{\text{precision}}$$

Combining these individual time complexities, we can determine the cumulative time complexity of the traversal distance algorithm as follows:

$$\begin{aligned} T &\in T_B \times (T_D + T_P) \\ &\in T_B \times O((|V_1| + |E_1|) \times (|V_2| + |E_2|)) + O((|V_1| \times |E_2|) + (|V_2| \times |E_1|)) \\ &\in T_B \times O(|V_1||E_2| + |V_2||E_1| + |V_1||E_1| + |V_2||E_2|) \\ &\in O(\log_2 \Delta) \times O(|V_1||E_2| + |V_2||E_1| + |V_1||E_1| + |V_2||E_2|) \\ &\in O(\log_2 \Delta \times (|V_1||E_2| + |V_2||E_1| + |V_1||E_1| + |V_2||E_2|)) \end{aligned}$$

This analysis establishes that the algorithm runs in $O(\log_2 \Delta \times (|V_1||E_2| + |V_2||E_1| + |V_1||E_1| + |V_2||E_2|))$ time [15]. Another property to revisit, as previously mentioned, is the traversal distance as an asymmetric measure. For the remainder of this thesis, it is asserted that [1]:

$$\forall (G_1 = G_2) \quad \delta_T(G_1, G_2) = \delta_T(G_2, G_1) = 0$$

2.6 Geometric Graph Edit Distance

A second method for measuring the distance between two geometric graphs involves calculating the edit distance between them. The edit distance between two objects is defined as the minimum number of operations required to transform one object into the other, where the operations include insertions, deletions, and relabeling [9]. Similar to how the Levenshtein distance measures the edit distance between two strings [10], the edit distance between two geometric graphs is referred to as the graph edit distance (GED).

For geometric graphs G_1 and G_2 , GED searches for a sequence of operations p that transforms $G_1 \rightarrow G_2$ such that the transformed graph $G'_1 = G_2$. An operation o may include deleting isolated vertices, inserting vertices, adding edges between existing vertices, deleting edges, and translating a vertex from one point to another. Each operation o has a corresponding cost function $Cost(o)$ for executing the operation. The cost for a sequence of operations is given by:

$$Cost(p) = \sum_{o_i \in p} Cost(o_i) \quad \text{where } p := \text{Set of Sequential Operations}$$

As a result, GED is defined as the cost of the least expensive path that transforms $G_1 \rightarrow G_2$.

$$GED(G_1, G_2) = \inf_{p \in P(G_1, G_2)} Cost(p)$$

$$P(G_1, G_2) := \text{Set of All } p \text{ that Transform } G_1 \rightarrow G_2$$

An advantage of GED over the traversal distance is its symmetry, which exists since the transformation of geometric graphs is reversible. One significant disadvantage is that computing GED is NP-hard, meaning that the time required to compute GED increases exponentially with the size of the geometric graphs [9].

3 Visualizing the Traversal Distance

This chapter builds on the traversal distance by demonstrating two visualization techniques. It starts with the free-space diagram for the weak Fréchet distance, explaining how it facilitates readers understanding of the concept. Additionally, a new visualization method for the traversal distance is introduced, addressing the challenges posed by the idea of a traversal distance free-space diagram.

3.1 Weak Fréchet Distance Free-Space Diagram

An important tool for visualizing the weak Fréchet distance is the FS_ϵ diagram. This visualization plays a role in enhancing the understanding of what the weak Fréchet distance algorithm computes. It also allows for the verification of whether cell boundaries within the diagram are being computed correctly [1]. The following discussion delves into the construction and interpretation of the FS_ϵ diagram, and how it builds on the visualization of the traversal distance free-space. The visualization in this section were generated using the Fréchet distance Python program documented in the appendix.

In Figures 8 and 9, the FS_ϵ diagram is demonstrated through a simple example involving two curves, C_1 and C_2 . C_1 is comprised of three line segments, while C_2 is comprised of four line segments. C_1 is aligned along the horizontal axis and C_2 along the vertical axis, with the free-space represented in white.

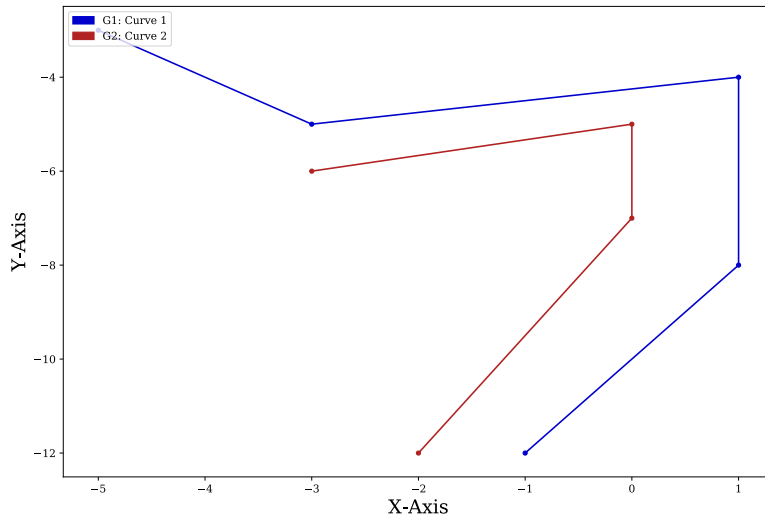


Figure 8: Curves C_1 and C_2 [15].

Assigning an arbitrary epsilon value, $\epsilon = 2$, to the FS_2 diagram, the resulting

visualization of the diagram can be observed as follows.

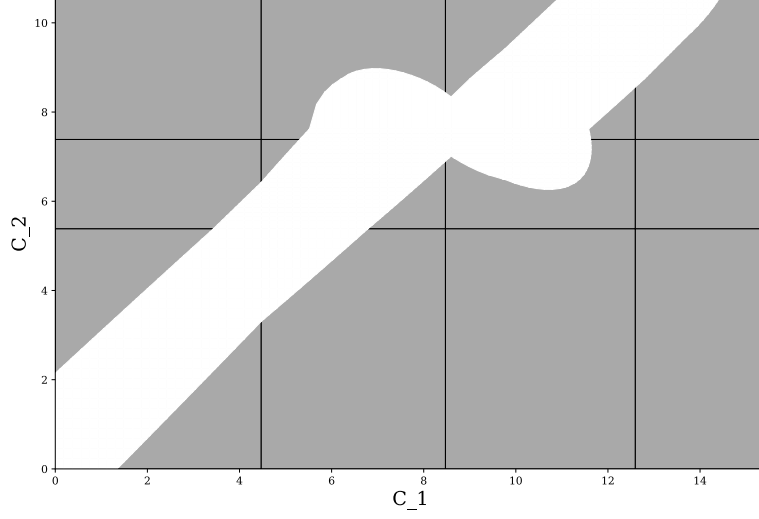


Figure 9: Figure 6 free-space diagram where $\epsilon = 2$ [15].

Observing Figure 8, the diagram consists of 12 free-space cells, matching the product of line segments in each curve: four in C_1 and three in C_2 . This epsilon value is not considered a weak Fréchet distance, however, since FS_2 does not cover every line segment inside C_1 .

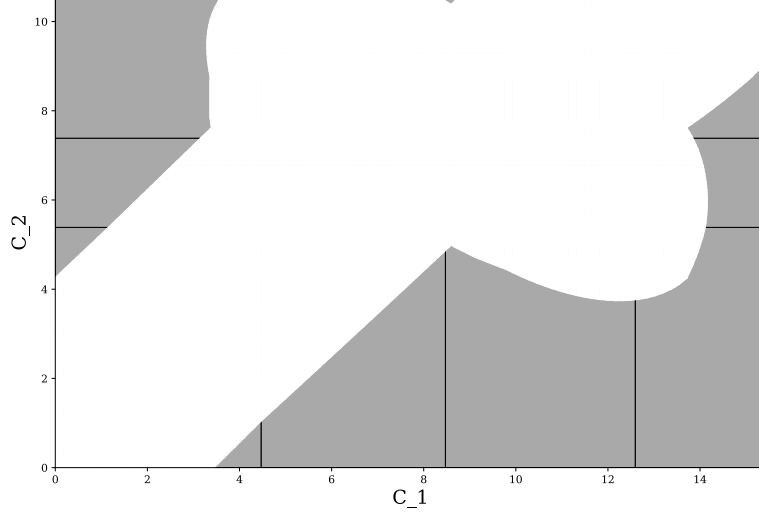


Figure 10: Free-space diagram where $\epsilon = 4$ [15].

Increasing epsilon to $\epsilon = 4$ in Figure 9 allows FS_4 to cover every line segment inside C_1 and C_2 , indicating that $\delta_F(C_1, C_2) \leq 4$ [15].

3.2 Traversal Distance Free-Space Visualization

Shifting focus to the traversal distance, a problem arises when visualizing the FS_ϵ diagram. The FS_ϵ diagram for the weak Fréchet distance can be effectively displayed in the Euclidean plane; this is achieved by mapping C_1 and C_2 along the X and Y axes in the Euclidean plane, respectively. However, this approach encounters a limitation when applied to geometric graphs. Such graphs cannot be similarly reduced to an axis in the Euclidean plane without causing overlapping of information. Consequently, this presents a unique challenge and opportunity when visualizing free-space for the traversal distance.

This section introduces a method for visualizing the traversal distance's free-space. Instead of displaying a free-space diagram with cells, this method focuses on visualizing the area specifically within the free-space cells.

Area of Free-Space Within a Cell In the case of two edges e_i and e_j , each from geometric graphs G_1 and G_2 , these two edges constitute a single free-space cell, labeled $FS_{\epsilon,i,j}$. Let A be the function that represents the area of the free-space cell $FS_{\epsilon,i,j}$, denoted by $A(FS_{\epsilon,i,j})$. Recall that $FS_{\epsilon,i,j}$ is represented as a square cell, constructed with four boundary walls, and the length from **start**

to end of a cell boundary falls within the range $[0, 1]$. Since a free-space cell is a square of length one, the range of $A(FS_{\epsilon,i,j})$ is bound by $[0, 1] \times [0, 1]$.

$$A(FS_{\epsilon,i,j}) \in [0, 1] \times [0, 1]$$

Since A is a values between $[0, 1]$, it can be interpreted as the percentage of area within a cell covered by free-space. For instance, if $A(FS_{\epsilon,i,j}) = 0.345$, this indicates that 34.5% of $FS_{\epsilon,i,j}$ is covered by free-space. $A(FS_{\epsilon,i,j}) = 0.0$ would indicate $FS_{\epsilon,i,j}$ is completely empty of free-space and $A(FS_{\epsilon,i,j}) = 1.0$ would indicate it is completely full.

Visualizing Free-Space Area To visualize the free-space area within $FS_{\epsilon,i,j}$, color in the quadrilateral area that lies between both edges e_i and e_j on their Euclidean plane. Consider this quadrilateral area the spatial relationship between these two edges. Figure 11 demonstrates the spatial relationship between an example pair of edges e_i and e_j .

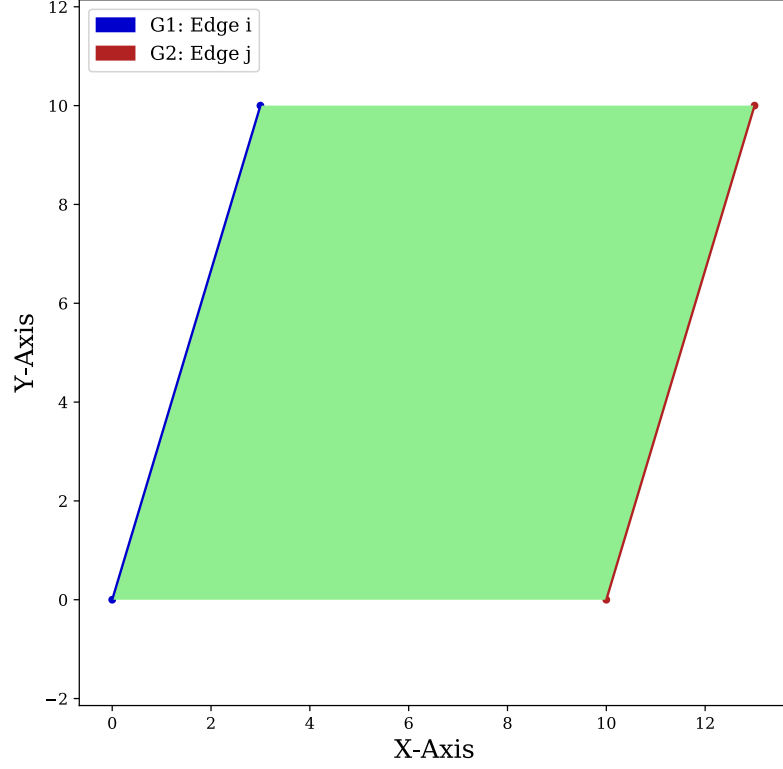


Figure 11: Highlighted space between e_{1_i} and e_{2_j} [15].

Furthermore, a transparency function α is applied to the color of the quadrilateral area. Such that the degree of transparency directly corresponds to the value of $A(FS_{\epsilon,i,j})$. This means that α visually represents the proportion of free-space covering cells. Figures 12 through 15 illustrate how, as ϵ increases from an empty free-space cell when $\epsilon = 0$ to a full free-space cell when $\epsilon = 15$, the color's transparency decreases.

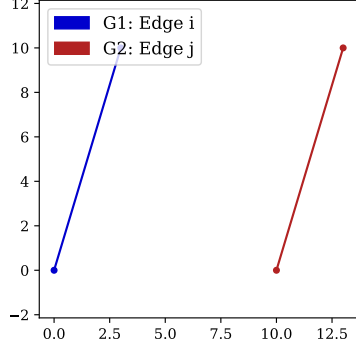


Figure 12: $\epsilon = 0$ [15].

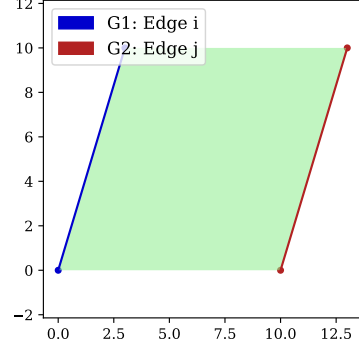


Figure 13: $\epsilon = 10.01$ [15].

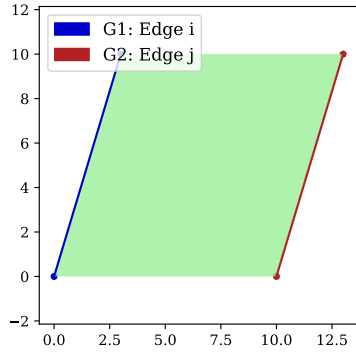


Figure 14: $\epsilon = 10.5$ [15].

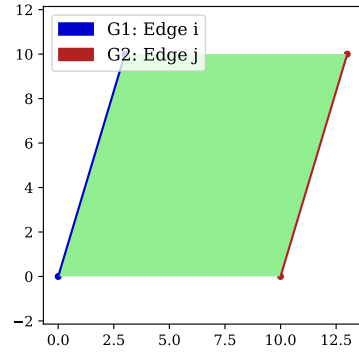


Figure 15: $\epsilon = 15$ [15].

This value α is amplified for overlapping areas. Consider a set of overlapping free-space areas $S = \{A_1, A_2, \dots, A_n\}$. The value of α for the intersection of these overlapping areas is calculated as follows [15]:

$$\alpha = 1 - \prod_{i=1}^n (1 - S_i)$$

3.3 Example of Traversal Distance Visualization

To demonstrate the overlapping property, consider an example involving a pair of geometric graphs representing more complex structures. This example involves a comparison between two distinct species of plant leaves, with each leaf constructed as a geometric graph. In these graphs, the edges represent both the outline of the leaf and its vein structure.

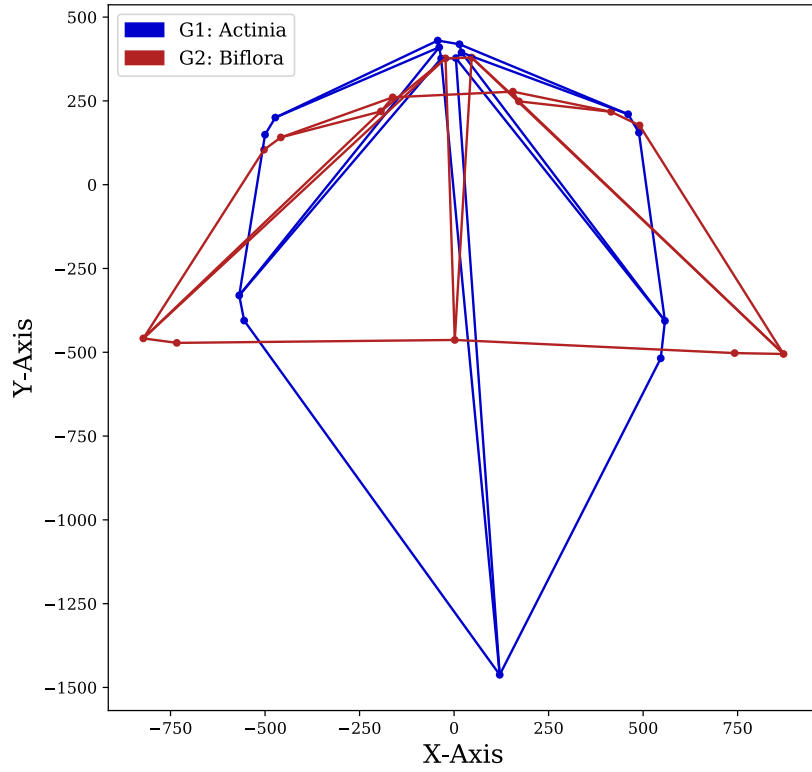


Figure 16: Two plant leaves, from the species *Actinia* and *Biflora*, as geometric graphs [13].

Having plotted the pair of geometric graphs, the free-space between both geometric graphs is now colored in for several values of ϵ .

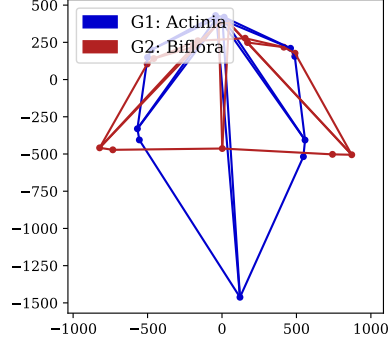


Figure 17: $\epsilon = 0$ [13].

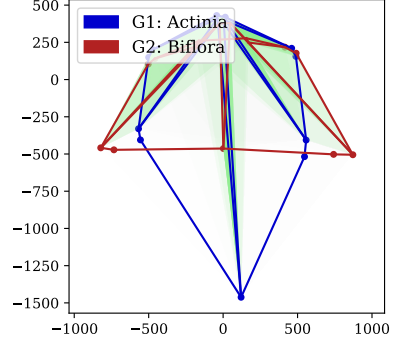


Figure 18: $\epsilon = 150$ [13].

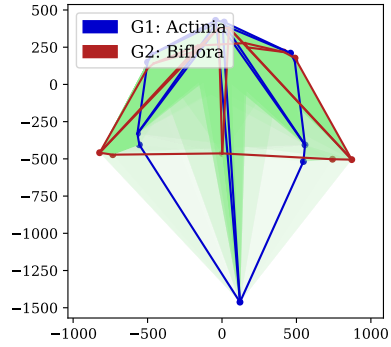


Figure 19: $\epsilon = 300$ [13].

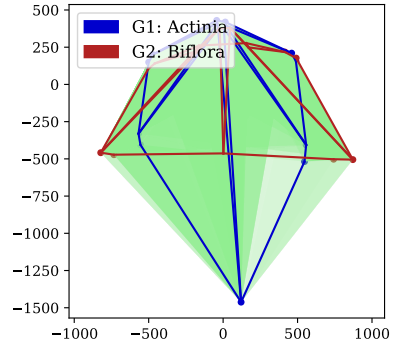


Figure 20: $\epsilon = 450$ [13].

Observing the visualizations, it becomes evident that the highlighted free-space expands as the value of epsilon increases [13].

4 Distance Measurements in Machine Learning

This chapter switches focus to concepts in machine learning: Euclidean distance, supervised learning, and classification. These topics will be important for differentiating different types of models within the field of machine learning. Followed by the application of distance measures in machine learning, with a particular focus on the k-nearest neighbors (k-NN) model. This segment of the chapter examines how distance metrics, such as the Euclidean distance, directly influence the behavior and performance of k-NN. The chapter introduces the concept of generalizing distance, explaining how k-NN, as defined here, performs when used with the traversal distance in the final chapter.

4.1 Euclidean Distance Between Two Points

The Euclidean distance is a conceptually straightforward distance metric in machine learning, tying to the earlier discussion about the properties of geometric graphs in the Euclidean plane. In the context of this thesis the Euclidean distance measures the straight-line distance between points in a two-dimensional space.

Definition 4. *Consider two points, $P_1 = (x_1, y_1)$ and $P_2 = (x_2, y_2)$ in a two-dimensional Euclidean plane. We define the Euclidean distance d as the shortest straight-line distance between two points in the plane, given by [3]:*

$$\|P_1 - P_2\| = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

4.2 Introducing the K-Nearest Neighbors Model

k-NN is a classification model in machine learning, specifically within supervised learning. As a non-parametric model, k-NN distinguishes itself by assuming the data does not have any specific underlying statistical distribution. This characteristic of being assumption-free renders k-NN not only basic but also an essential model for grasping the concepts in machine learning classification.

Classification in Supervised Learning Supervised learning in machine learning refers to the process where a model is trained using a labeled dataset. Where the term supervised implies that the model learns from the input data X and output label y . In other words, the supervision of learning can be thought of as a mapping function $model(X) = y$.

Classification in supervised learning involves categorizing data into predefined classes. A model predicts the class of a new observation x_{new} by analyzing X , then assigns x_{new} a predicted label \hat{y} .

Classification can be separated into two formats: binary and multi-class. In binary classification, the model classifies observations into one of two classes. While multi-class classification involves categorizing observations into one of multiple classes. This thesis focuses only on multi-class classification [6].

K-Nearest Neighbors Algorithm The process of predicting a class, for an observation x_{new} , can be broken down into the following steps:

1. Choose the Number of Neighbors: Select the number of nearest neighbors, k , which will influence the prediction. How to determine a value for k is discussed in the appendix.”
2. Calculate Distances: For the observation x_{new} , compute the Euclidean distance between x_{new} and each observation in X .
3. Identify Nearest Neighbors: Sort all calculated distances in ascending order. Then select the top k nearest observations from the dataset.
5. Predict the Classification: Determine the predicted label \hat{y} for the observation x_{new} by taking the most common label among the k neighbors. This step is known as a majority voting ensemble in machine learning.

The formal definition of the k-NN algorithm is presented as Algorithm 3 in [12].

Algorithm 3 KNearestClassifier

Require: Dataset X

Require: Labels y

Require: New observation x_{new}

Initialize number of neighbors k

Initialize distance d

Initialize list *distances*

Begin Algorithm

for each point x_i in the training data X **do**

 Calculate distance $d(x_i, x_{new})$

 Append $d(x_i, x_{new})$ to *distances*

end for

Sort *distances* in ascending order

Select top k nearest neighbors from *distances*

Aggregate labels of selected neighbors

$\hat{y} :=$ Most common label

return \hat{y}

End Algorithm

A sixth step of the algorithm, not previously discussed, addresses situations where majority voting leads to a tie. This occurs when two or more labels are equally most common among the aggregated labels [6].

Time Complexity Let the time complexity of our distance algorithm to be $O(1)$. Assuming the sorting algorithm used is Quick Sort, the time complexity of the sorting process can be defined as $O(|X| \log_2 |X|)$ [4], where $|X|$ represents the number of observations in our training data. Therefore, the overall time

complexity of a k-NN prediction can be expressed as follows [6]:

$$T_{KNN} \in O(|X| + |X| \log_2 |X|)$$

4.3 Evaluating K-Nearest Neighbors Predictions

This section focuses on two tools used when assessing the performance of classification models: the train-test dataset split and the evaluation metrics of precision and recall. Further describing the methodology behind partitioning datasets into training and testing subsets, and subsequently explore how precision and recall scores effectively measure a model’s predictive accuracy [6]. Additionally, the practical application of these evaluation techniques will be demonstrated using the scikit-learn Python library [12].

Training and Testing Datasets To assess the effectiveness of k-NN in predicting new observations, datasets are divided into two disjoint sets: a training dataset and a testing dataset. This division is uniformly random to establish unbiased distribution. The test dataset serves as independent data to evaluate the accuracy of the model’s predictions, specifically by measuring how many predicted values (\hat{y}) are correctly classified as the actual values (y) [6]. In scikit-learn, the function used to perform this dataset split is structured as follows [12]:

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2)
```

Precision and Recall of Predictions In classification, precision and recall are metrics used to evaluate the accuracy and relevance for a set of model predictions. Precision, in this context, measures the proportion of correctly predicted positive observations to the total predicted positive observations. It is a key indicator of a model’s ability to minimize false positives. Recall, on the other hand, assesses the proportion of actual positive observations that were correctly predicted by the model, thus reflecting its capability to minimize false negatives. For multi-classification problem, the Precision and Recall of a set of predictions is defined as [5]:

Definition 5. *Given a set of multiple classes $C = \{c_i | i \in \mathbb{N}\}$, within a set of predictions.*

- **Precision** for a specific class c_i is defined as the ratio of the number of true positive instances TP_i to the total number of instances predicted as belonging to that class, which is the sum of the true positives and false positives FP_i :

$$P_i = \frac{TP_i}{TP_i + FP_i}$$

- **Recall** for class c_i is the ratio of the number of true positive instances TP_i to the actual number of instances of that class in the data, which is the sum of the true positives and false negatives FN_i :

$$R_i = \frac{TP_i}{TP_i + FN_i}$$

When using the `y_test` dataset and considering `y_pred` as the set of model predictions \hat{y} , the calculation of precision and recall using scikit-learn can be computed as follows [12]:

```
from sklearn.metrics import precision_score, recall_score

precision = precision_score(y_test, y_pred, average='macro')
recall = recall_score(y_test, y_pred, average='macro')
```

5 Applying the Traversal Distance to Classification Problems

To test the precision of the traversal distance within the framework of the k-nearest neighbors (k-NN) algorithm, the inherent asymmetry of the traversal distance must be addressed. Therefore, a symmetric adaptation suitable for k-NN is proposed. This adjustment enables the k-NN algorithm to classify a dataset of English letters represented as geometric graphs using traversal distance [14]. The goal is to evaluate the hypothesis that traversal distance can precisely classify geometric graphs. This is determined by comparing the performance of the traversal distance k-NN model with that of the graph edit distance (GED) k-NN model.

5.1 Symmetric Case of the Traversal Distance

To incorporate the traversal distance into the k-NN algorithm, it is necessary to first address the asymmetry of the traversal distance. Meaning, the distance from G_1 to G_2 may not equal the distance from G_2 to G_1 , resulting in two distinct distances. k-NN operates under the assumption that the distance metric implemented is symmetric [6]. Consequently, a symmetric variant of the traversal distance must be defined.

To develop a symmetric distance metric for k-NN, a function must be devised that combines the two asymmetric distances produced by the traversal distance into a single distance measure [1]. The design of a function should be tailored to the specific requirements of a dataset [6]. Given that this chapter concentrates on the comparison of English characters, it operates under the presumption that the geometric graphs being compared have equivalent magnitudes [14].

For the distance metric of k-NN, the function will be defined as the maximum of the two distances produced by the traversal distance. Taking the maximum value ensures that the distance, denoted by ϵ , completely covers both geometric graphs during `projection_check`.

Theorem 1. *Let the symmetric traversal distance between two geometric graphs G_1 and G_2 be defined by the equation:*

$$\delta_{ST}(G_1, G_2) = \max\{\delta_T(G_1, G_2), \delta_T(G_2, G_1)\}$$

*Then, it holds that the distance metric ϵ associated with $\delta_{ST}(G_1, G_2)$ passes the **projection_check** for both $\delta_T(G_1, G_2)$ and $\delta_T(G_2, G_1)$, ensuring that ϵ fully covers both G_1 and G_2 for the symmetric case.*

Proof. For the symmetric traversal distance between two geometric graphs G_1 and G_2 , assume without loss of generality:

$$\delta_{ST}(G_1, G_2) = \delta_T(G_1, G_2) \geq \delta_T(G_2, G_1)$$

Where $\delta_T(G_1, G_2) = \epsilon_1$ and $\delta_T(G_2, G_1) = \epsilon_2$. This implies that $\epsilon_1 \geq \epsilon_2$. By **Definition 3**, $\delta_T(G_1, G_2)$ traverses entirely over G_1 and traverses partially

over G_2 . Given that the free-space is a monotone function of ϵ , this statement holds for any $\epsilon \in [\epsilon_2, \infty)$. If $\epsilon_1 \geq \epsilon_2$, then $\epsilon_1 \in [\epsilon_2, \infty)$, ensuring $\delta_T(G_1, G_2)$ traverses over the entirety of both G_1 and G_2 . Therefore, $\delta_{ST}(G_1, G_2)$ covers both graphs for the purpose of **projection_check**. □

5.2 K-Nearest Neighbors Using the Traversal Distance

Having established a symmetric traversal distance, it is now possible to predict classifications using the English letter dataset. In the context of the English Letter dataset, X comprises of, English letters, stored as geometric graphs rather than traditional points, y denotes the class of letter, and d serves as the metric for quantifying the distance between any two geometric graphs.

The dataset comprises 2,250 labeled geometric graphs, each representing a distorted drawing of an English letter. These drawings are categorized into 15 distinct classes of 150 observations. Each class corresponds to one letter such that $C = \{A, E, F, H, I, K, L, M, N, T, V, W, X, Y, Z\}$ [14].

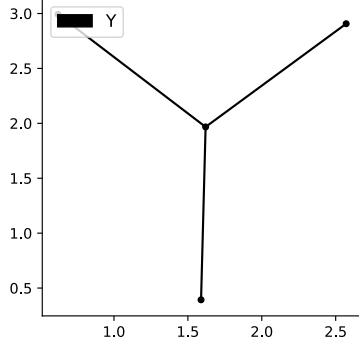


Figure 21: Class Y [14].

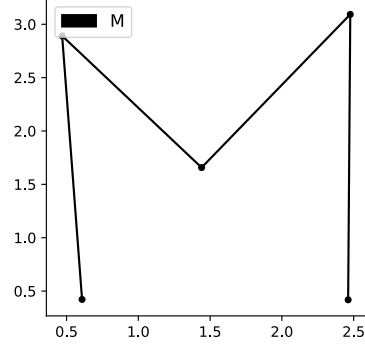


Figure 22: Class M [14].

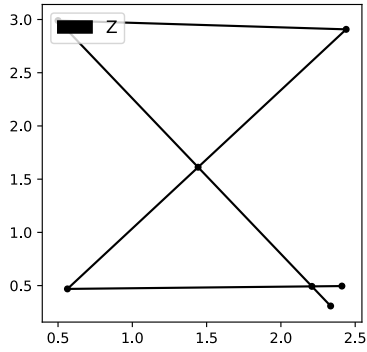


Figure 23: Class Z [14].

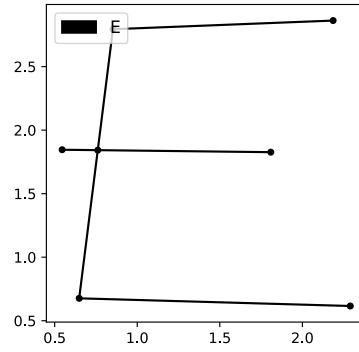


Figure 24: Class E [14].

The dataset is partitioned into independent training and testing datasets using a uniform random distribution, allocating approximately 80 percent of the observations to the training set and 20 percent to the testing set [6]. Leaving the distribution of observation across classes, training and testing as follows [14]:

Observations in Training Dataset: 1794

Observations in Testing Dataset: 449

For k-NN, the parameter specifying the number of neighbors, is set to $k = 7$. Meanwhile, for the traversal distance calculation, the `binarySearch` function's parameters are configured as follows: *left* = 0, *right* = 3, and *precision* = 0.001.

```
model = KNeighborsClassifier(n_neighbors=7, mean='max',
                             left=0, right=3, precision=0.001)
```

The model is first fitted on the training dataset, then predicts the observations in the test dataset, yielding 449 predicted classes from the test dataset, denoted as \hat{y} [6]. To evaluate the precision and recall of the $|C| = 15$ classes between \hat{y} and the true labels y , both Macro-Average Precision and Macro-Average Recall are calculated. These metrics are computed by first calculating the precision for each class independently then taking the average of these precision scores. Use of this method allows for datasets with class imbalances while ensuring that errors across all classes are treated equally [5].

$$\text{Macro-Average Precision} = \frac{1}{|C|} \sum_{i=1}^{|C|} P_i \quad \text{Macro-Average Recall} = \frac{1}{|C|} \sum_{i=1}^{|C|} R_i \quad (1)$$

After running the test for approximately 7 hours, the results showed a Macro-Average Precision of 89.5 percent and a Macro-Average Recall of 87.6 percent. These results resemble those obtained when using k-NN with the graph edit distance [14].

5.3 K-Nearest Neighbors Using Graph Edit Distance

To compare the results obtained using traversal distance, it is useful to examine a benchmark test conducted with GED. In their seminal work, Kaspar Riesen and Horst Bunke investigate the efficacy of the k-NN model, employing GED, in classifying geometric graph representations of the English alphabet. This comparison serves as the reference point for assessing the performance of traversal distance.

In their investigation, Riesen and Bunke utilize a dataset comprising 6,750 observations of English letters, with 750 observations uniformly distributed across 15 letter classes. To evaluate the performance of the GED k-NN model, the dataset was divided into equal parts for validation, training, and testing. As a result, the model achieved a macro-average precision of 99.6 percent on the test dataset, showcasing the high performance of GED in this context [14].

When compared, the GED k-NN model outperforms the traversal distance k-NN model by 11.3 percent.

6 Conclusion

This thesis examined the application of geometric graph distances within computer science, with a specific focus on the traversal distance and the algorithm for computing it. The aim was to assess the utility of the traversal distance in classifying geometric graphs within the English letter dataset, a challenge in supervised machine learning. Comparative analysis revealed that the k-NN model based on traversal distance achieved precision comparable to that of the GED-based k-NN model by Riesen and Bunke, despite the traversal distance being polynomial and GED being NP-Hard. This result indicates that the traversal distance is an effective method for classifying geometric graphs in supervised machine learning.

The research presented several key findings. The traversal distance was defined by incorporating elements of geometric graph theory, weak Fréchet distance, and free-space diagrams. The algorithm's steps were detailed, and its time complexity was determined to be polynomial. Additionally, the thesis proposed a method for visualizing free-space areas related to the traversal distance, demonstrating its effectiveness with examples from the plant leaf dataset. A symmetric version of the traversal distance was defined, wherein taking the maximum value of both asymmetric distances ensures that the free-space entirely covers both geometric graphs. The efficacy of the traversal distance in enhancing the k-NN model's performance for classifying geometric graphs, as evidenced in the English letter dataset, was measured with precision and recall metrics.

For future work, the thesis suggests avenues for refining the traversal distance. The symmetric traversal distance algorithm could be optimized by preventing the recalculation of free spaces already computed within the maximum function. Moreover, there is potential for enhancing the overall runtime efficiency of the traversal distance algorithm through the adoption of a lower-level programming language. In conclusion, the thesis underscores that distances in geometric graph theory represent a relatively new field within mathematics that is currently evolving. The application of these distances, particularly in computer science for digital road maps and increasingly within supervised machine learning, is expected to broaden as the efficiency of these algorithms improves and their comparative advantages are tested against existing software and models.

Appendix

This appendix documents the Python packages and Jupyter Notebooks written to support this thesis. All supporting materials are publicly available in the `compgeomTU/will_rodman_thesis` GitHub repository.

FSDVis Python Package

The FSDVis Python Package computes then visualizes the free-space diagram for the weak Fréchet distance. The package is located at `compgeomTU/FSDVis`. Using Figure 4 as an example, the following script from Jupyter Notebook `will_rodman_thesis/images/image_generator.ipynb` imports the package and plots a free-space diagram.

```
import matplotlib.pyplot as plt
from FSDVis.Curve import Curve
from FSDVis.FreeSpaceDiagram import FreeSpaceDiagram

filepath_1 = 'examples/pair/aside'
filepath_2 = 'examples/pair/bside'
curve_1 = Curve(filepath_1)
curve_2 = Curve(filepath_2)
epsilon = 6

fsd = FreeSpaceDiagram(curve_1, curve_2, n_approximation=25)
fsd.buildFreeSpace(epsilon)
fsd.buildCells()

fig, ax = fsd.plotFreeSpace()
fig.set_size_inches(4, 4)
ax.set_xlabel("Edge i", fontsize=axis_fontsize_small_square)
ax.set_ylabel("Edge j", fontsize=axis_fontsize_small_square)
plt.show()
```

TraversalDistance Python Package

The TraversalDistance Python Package contains the traversal distance algorithm, traversal distance visualizer and k-NN algorithm implementing the symmetric traversal distance. The package located at `compgeomTU/TraversalDistance` contains the following files.

Using Figure 18 as an example, the following script from Juniper Notebook `will_rodman_thesis/images/image_generator.ipynb` imports the package and plots the traversal distance free-space area.

```
import matplotlib.pyplot as plt
from TraversalDistance.Graph import Graph
from TraversalDistance.Visualize import Visualize
```

File	Description
BinarySearch.py	BinarySearch class function of the traversal distance algorithm.
CalFreeSpace.py	Function for computing then returning the $[start, end]$ boundary for free-space cell walls.
FreeSpaceGraph.py	DFSTraversalDistance class function of the traversal distance algorithm.
Graph.py	Class data structure for storing geometric graphs.
KNeighborsClassifier.py	Custom k-NN class that implements the symmetric traversal distance.
LineIntersection.py	Supporting class for the projection_check function.
Visualize.py	Class responsible for visualizing the traversal distance free-space area.

Table 1: TraversalDistance Python Package Files.

```

filepath_1 = 'examples/plant/Pact1_actinia_3'
filepath_2 = 'examples/plant/Pbif1_biflora_1'
graph_1 = Graph(filename=filepath_1, name='Actinia')
graph_2 = Graph(filename=filepath_2, name='Biflora')

visualize = Visualize(graph_1, graph_2, epsilon=epsilon)

fig, ax = visualize.plot_freespace(legend_fontsize=legend_fontsize)
fig.set_size_inches(4, 4)
plt.show()

```

Traversal Distance K-NN Model Test: Jupyter Notebooks

The analysis of the traversal distance k-NN model, conducted on the English letter dataset, was executed within the Jupyter notebook located at `will_rodman_thesis/letter_data/letter_knn`. The model required approximately seven hours to predict all the values in the test dataset. The predicted and actual classes derived from this test were logged and subsequently analyzed in the Jupyter notebook `will_rodman_thesis/letter_data/letter_knn_analysis`. This analysis, illustrated in Figure 23 demonstrates how precision and recall metrics vary in relation to the value of k .

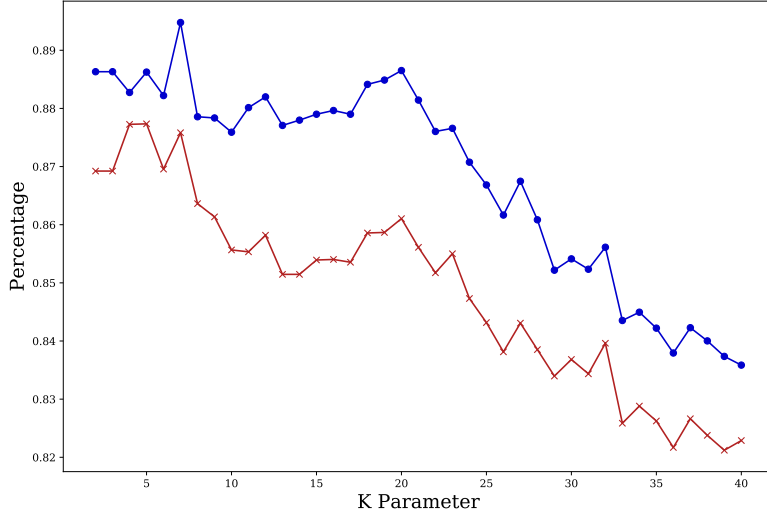


Figure 25: Traversal distance k-NN model precision and recall scores on the English letter dataset

Upon reviewing Figure 23, both precision and recall exhibit a consistent decline as k increases from 20 to 40, with peak precision occurring at $k = 7$.

References

- [1] Helmut Alt, Alon Efrat, Günter Rote, and Carola Wenk. Matching planar maps. *Journal of Algorithms*, 49(2):262–283, 2003.
- [2] Daniel Chen, Christian Sommer, and Daniel Wolleb. Fast map matching with vertex-monotone fréchet distance. In *Algorithmic Approaches for Transportation Modeling, Optimization, and Systems*, 2021.
- [3] Ivan Dokmanic, Reza Parhizkar, Juri Ranieri, and Martin Vetterli. Euclidean distance matrices: Essential theory, algorithms, and applications. *IEEE Signal Processing Magazine*, 32(6):12–30, November 2015.
- [4] Shalosh B. Ekhad and Doron Zeilberger. A detailed analysis of quicksort running time, 2019.
- [5] Margherita Grandini, Enrico Bagli, and Giorgio Visani. Metrics for multi-class classification: an overview, 2020.
- [6] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2 edition, 2009.
- [7] Erfan Hosseini. Graphsamplingtoolkit: Compare roadmaps or evaluate reconstructed maps with graph sampling toolkit. <https://github.com/Erfanh1995/GraphSamplingToolkit>, 2021.
- [8] Kunal Kumar, Tushar Kumar, and Gargi Chakraborty. Decomposed algorithm for reducing time complexity in binary search, 04 2021.
- [9] Sushovan Majhi and Carola Wenk. Distance measures for geometric graphs, 2022.
- [10] Evgenii Ofitserov, Vasily Tsvetkov, and Vadim Nazarov. Soft edit distance for differentiable comparison of symbolic sequences, 2019.
- [11] János Pach. 10 geometric graph theory. 2017.
- [12] Pedregosa, F. and Varoquaux, G. and Gramfort, A. and Michel, V. and Thirion, B. and Grisel, O. and Blondel, M. and Prettenhofer, P. and Weiss, R. and Dubourg, V. and Vanderplas, J. and Passos, A. and Cournapeau, D. and Brucher, M. and Perrot, M. and Duchesnay, E. *Scikit-learn: Machine Learning in Python*. Scikit-learn developers, 2023.
- [13] Sarah Percival, Joyce Onyenedum, Daniel Chitwood, and Aman Husbands. Topological data analysis reveals core heteroblastic and ontogenetic programs embedded in leaves of grapevine (vitaceae) and maracuyá (passifloraceae). *PLoS computational biology*, 20:e1011845, 02 2024.
- [14] Kaspar Riesen and Horst Bunke. Iam graph database repository for graph based pattern recognition and machine learning. 2008.

- [15] Will Rodman. `will_rodman_thesis`: Traversal distance python library. https://github.com/compgeomTU/will_rodman_thesis, 2024.
- [16] Carola Wenk. Weak fréchet code. <https://www.cs.tulane.edu/~carola/research/code.html>, 2018.