# CS 533 Project 3

Freddie Awuah-Gyasi      Jaime Gould      Lasair Servilla      Qinghong Shao

`fawuahgyasi5@unm.edu`      `egould@unm.edu`      `lservilla@unm.edu`      `qinghongshao@unm.edu`

due: April $21^{st}$ 2025

## 1 Functions and Pseudocode

When we initially began to implement the Hopkins Statistic in **Project 2**, we noticed that the algorithm could easily be inefficient when searching for the nearest neighbor in a dataset to a given point:

---
**Algorithm 1** Hopkins Statistic

---
1: **function** HOPKINSSTAT($D, m$)
2:      $D_m \leftarrow$ RANDOMSUBSET($D$, $m$)          $\triangleright \mathcal{O}(m)$
3:      $R_m \leftarrow$ RANDOMPOINTS($m$, $x_{range}$, $y_{range}$)          $\triangleright \mathcal{O}(m)$
4:      $r \leftarrow \{ \ |R_i -$ NEARESTNEIGHBOR($R_i$, $D$)$| \ \} \ \forall \ R_i \in R_m$      $\triangleright \mathcal{O}(m) \times \mathcal{O}(\text{NEARESTNEIGHBOR})$
5:      $p \leftarrow \{ \ |D_i -$ NEARESTNEIGHBOR($D_i$, $D$)$| \ \} \ \forall \ D_i \in D_m$      $\triangleright \mathcal{O}(m) \times \mathcal{O}(\text{NEARESTNEIGHBOR})$
6:      $r_{\text{sum}}^2 \leftarrow \sum_{i=1}^{m} r_i^2 \ \ \forall \ r_i \in r$          $\triangleright \mathcal{O}(m)$
7:      $p_{\text{sum}}^2 \leftarrow \sum_{i=1}^{m} p_i^2 \ \ \forall \ p_i \in p$          $\triangleright \mathcal{O}(m)$
8:      $H \leftarrow \frac{r_{\text{sum}}^2}{p_{sum}^2 + r_{\text{sum}}^2}$          $\triangleright \mathcal{O}(1)$
9:      **return** $H$
10: **end function**

---

If we simply iterate over all points in a given dataset $D$ (where the size of the dataset $|D| = n$) for each point $r \in R_m$ and $p \in D_m$ ($|R_m| = |D_m| = m$) in order to determine which point is closest, we end up with a complexity of $\mathcal{O}(m) \times \mathcal{O}(\text{NEARESTNEIGHBOR}) = \mathcal{O}(m \times n)$ for a single pass of the algorithm alone. When we then compute the Hopkins Statistic for a bootstrap size of $B$, the complexity of this algorithm increases to $\mathcal{O}(B \times m \times n)$. This can be **very** slow depending on the size of $B$, $n$, and $m$.

In order to improve efficiency of our algorithm we created two classes, BOUNDINGBOX and GRID. BOUNDINGBOX merely contains the maximum and minimum of the $x$ and $y$ ranges over which the points can exist, while GRID organizes the points of a given dataset into a grid structure with bins of size *partition_size* (see Figure 1, left).

---
1: **class** BOUNDINGBOX
2:      $x_{min}$ : FLOAT          $\triangleright$ Our datasets have $x_{min} = -10000$
3:      $x_{max}$ : FLOAT          $\triangleright x_{max} = 10000$
4:      $y_{min}$ : FLOAT          $\triangleright y_{min} = -10000$
5:      $y_{max}$ : FLOAT          $\triangleright y_{max} = 10000$
6: **end class**

---

```
 1: class GRID
 2:     function INIT(points : LIST, bbox : BOUNDINGBOX, partition_size : FLOAT)
 3:         x_min ← bbox.x_min − partition_size                                    ▷ O(1)
 4:         y_min ← bbox.y_min − partition_size                                    ▷ O(1)
 5:         i_{x_max} ← FLOOR((bbox.x_max − x_min)/partition_size)                 ▷ O(1)
 6:         i_{y_max} ← FLOOR((bbox.y_max − y_min)/partition_size)                 ▷ O(1)
 7:         grid ← array of empty arrays, where ||grid|| = i_{x_max} × i_{y_max}   ▷ O(i_{x_max} × i_{y_max}) = O(1)
 8:         for p ∈ points do
 9:             x_index ← FLOOR((p.x − x_min)/partition_size)
10:             y_index ← FLOOR((p.y − y_min)/partition_size)
11:             grid[x_index][y_index] ← p
12:         end for                                                                ▷ O(n)
13:     end function
14:
15:     function GETNEIGHBORS(p, neighborhood)
16:         x_idx ← MAX(MIN(FLOOR((p.x − x_min)/partition_size), i_{x_max}), 0)    ▷ O(1)
17:         y_idx ← MAX(MIN(FLOOR((p.y − y_min)/partition_size), i_{y_max}), 0)    ▷ O(1)
18:         neighbors ← grid[x_idx][y_idx]                                         ▷ O(1)
19:
20:         for MAX(x_idx − neighborhood, 0) ≤ i ≤ MIN(x_idx + neighborhood, i_{x_max}) do
21:             for MAX(y_idx − neighborhood, 0) ≤ j ≤ MIN(y_idx + neighborhood, i_{y_max}) do
22:                 neighbors ← neighbors + grid[i][j]
23:             end for
24:         end for                                              ▷ ≤ O(neighborhood × neighborhood)
25:         return neighbors
26:     end function
27:
28:     function NEARESTNEIGHBOR(p)
29:         neighbors ← GETNEIGHBORS(p, 0)                                         ▷ O(1)
30:         if |neighbors| < 2 then neighbors ← GETNEIGHBORS(p, 1)                 ▷ O(1 × 1) = O(1)
31:         end if
32:         if |neighbors| < 2 then neighbors ← GETNEIGHBORS(p, 2)       ▷ O(2 × 2) = O(4) = O(1)
33:         end if
34:         if |neighbors| < 2 then neighbors ← GETNEIGHBORS(p, 3)       ▷ O(3 × 3) = O(9) = O(1)
35:         end if
36:         if |neighbors| < 2 then neighbors ← GETNEIGHBORS(p, 4)       ▷ O(4 × 4) = O(16) = O(1)
37:         end if
38:         if |neighbors| < 2 then neighbors ← points                            ▷ O(1)
39:         end if
40:
41:         nearest ← ∅
42:         dist ← ∞
43:         for v ∈ neighbors do
44:             new_dist ← EUCLIDEANDISTANCE(v, p)
45:             if new_dist < dist then
46:                 nearest ← v
47:                 dist ← new_dist
48:             end if
49:         end for                                                                ▷ O(|neighbors|)
50:         return nearest
51:     end function
52: end class
```

## 2   Time Complexity Analysis

GRID is initialized once for a given dataset $D$ ($|D| = n$) before the Hopkins Statistic is computed. Upon initialization, a grid data structure is created (GRID, line 7) and every point $p \in D$ is placed into it according to its $x$ and $y$ parameters (lines 8-12). Setting up GRID has initial time complexity $\mathcal{O}(n)$.
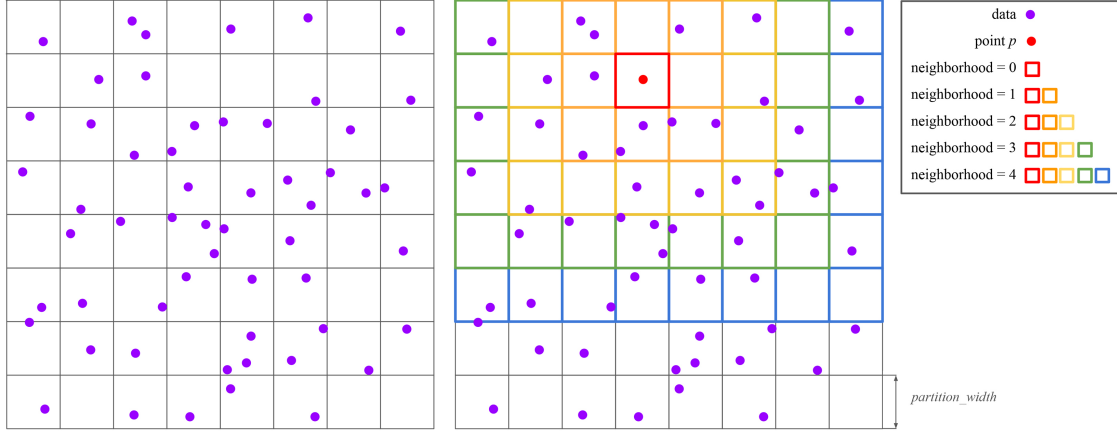


Figure 1: Left: an example of 2D data points in our GRID data structure. Right: the neighborhood of a given point $p$ (in red) at different neighborhood sizes.

Once GRID has been initialized, we then use its NEARESTNEIGHBOR function to find the neighbor nearest to a given point $p$ in our HOPKINSSTATISTIC algorithm. In the best case scenario of this function, there are points in the closest neighborhood of $p$. However, if there are less than 2 points[1] in the list of neighbors we retrieve, we search for the next largest neighborhood = 1 (see Figure 1, right), and so on up to a neighborhood of 4. If the number of points is still $|neighbors| < 2$, we set our neighbors to be equal to the list of total points in the dataset $D$. We then compare the euclidean distance of our point $p$ to each neighbor in $neighbors$, resulting in complexity of NEARESTNEIGHBOR depending on the size of $neighbors$, $\mathcal{O}(|neighbors|)$.

In the best case scenario, a smaller neighborhood yields a number of points $|neighbors| = 2$, giving complexity $\mathcal{O}(2) = \mathcal{O}(1)$. If the smaller neighborhood fails to retrieve at least two points or happens to contain all points in the dataset due to the way the data is distributed in $D$, the euclidean distance is computed for all points in $D$, resulting in $\mathcal{O}(n)$. Therefore, the complexity of our algorithm for a single run of the the Hopkins Statistic is:

$$\mathcal{O}(\text{HOPKINSSTATISTIC}) = \mathcal{O}(m) \times \mathcal{O}(\text{NEARESTNEIGHBOR}) = \mathcal{O}(m) \times \mathcal{O}(|neighbors|)$$

$$\mathcal{O}(m) \times \mathcal{O}(1) \leq \mathcal{O}(\text{HOPKINSSTATISTIC}) \leq \mathcal{O}(m) \times \mathcal{O}(n)$$
$$\mathcal{O}(m) \leq \mathcal{O}(\text{HOPKINSSTATISTIC}) \leq \mathcal{O}(m \times n)$$

The bootstrapping portion of our code, BOOTSTRAP, which computes the HOPKINSSTATISTIC $B$ times then has complexity $\mathcal{O}(\text{BOOTSTRAP}) = \mathcal{O}(B) \times \mathcal{O}(\text{HOPKINSSTATISTIC})$:

$$\mathcal{O}(B) \times \mathcal{O}(m) \leq \mathcal{O}(\text{BOOTSTRAP}) \leq \mathcal{O}(B) \times \mathcal{O}(m \times n)$$
$$\mathcal{O}(B \times m) \leq \mathcal{O}(\text{BOOTSTRAP}) \leq \mathcal{O}(B \times m \times n)$$

---

[1] $D_m \subset D$, and the nearest neighbor of a point $p \in D_m$ in $D$ should not be itself. Therefore we check for a set of neighbors of at least 2 so that $neighbors$ does not include only $p$.

# 3 Experiments for Testing Complexity

In order to determine whether our time complexity analysis is correct, we have designed experiments in which we alter each of our variables that we have hypothesized mainly impact the complexity: $B$, $m$, and $n$ ($n$ is fixed at 5000 for each dataset in **Project 2**. We will create random subsets of varying size $n$ from our original datasets in order to simulate different dataset sizes).

We utilize linear regression in order to test whether our hypotheses are correct. We can format our Cost function $C$ in the form of a linear equation:

$$r(x) = \beta_\circ + \beta_1 \cdot x$$

$$\beta_\circ + \beta_1(B \times m) \leq C(B, m, n) \leq \beta_\circ + \beta_1(B \times m \times n)$$

Where $\beta_\circ$ is the constant time the program takes when $B = m = n = 0$, $\beta_1$ is the rate at which the time increases as $B$, $m$, and $n$ increase (at high values of $B$, $m$, and $n$ when $\mathcal{O}(B \times m \times n)$ has the greatest impact), and $y$ is the cost function, or the total time the program takes to run.

## 3.a Hypotheses

**Null Hypothesis.** *($H_\circ$) There is no relationship between B, m, n and time; therefore $\beta_1 = 0$.*

**Alternative Hypothesis.** *($H_1$) There is a linear relationship between B, m, n and time; parameter $\beta_1$ is some constant value.*

We can use the $F$-Test to determine if our expected value $E(\beta_1) = 0$ under the null hypothesis falls within the 95% confidence interval of our computed $\beta_1$ value, where:

## 3.b Experiments

We will run our code while changing our parameter $B$ at intervals of 1000 ($1000 \leq B \leq 11000$), our parameter $m$ at intervals of 100 ($100 \leq m \leq 1100$), and our parameter $n$ at intervals of 300 ($2000 \leq n \leq 5000$) for each of our 4 datasets ($D_1, D_2, D_3, D_4$) from **Project 2**. This will result in $10 \times 10 \times 10 = 1000$ runs for each dataset.

We will record the time it takes for each individual run. We expect to see the time increasing as each of our parameters increases according to our time complexity analysis between lower bound $\mathcal{O}(B \times m)$ and upper bound $\mathcal{O}(B \times m \times n)$, so we will compute the linear regression fit and estimates for $\hat{\beta}_\circ$ and $\hat{\beta}_1$ and their 95% confidence intervals for each dataset and plot the time of our runs as a function of $x = B \times m \times n$.