

Heuristics for Master Thesis

Håvard Notland

February 2020

1 Introduction

In the modern era, many jobs which were never imagined 10 years ago are starting to appear. One of these jobs is to collect and charge electric scooters.

In several major cities across the world, you can now rent an electric scooter on almost every street corner, and it's popularity has exploded in a very short amount of time. Locating a scooter and renting is easily done via an app on your mobile phone. Once your scooter has taken you to your destination, you may simply leave it on the sidewalk, as long as it is in within a designated zone. Later, someone else may pick up the scooter and continue to use it. Once it has run out of power, it needs to be collected and charged.

This logistics task can be tackled in various ways: Some companies hire designated drivers that pick up the scooters and take them to one of several charging stations throughout the city. More commonly used is a crowd-sourcing effort where anyone can sign up to become a scooter collector, and get paid per scooter collected, charged and returned to the streets. These independent contractors are commonly referred to as "Bird-catchers" or "Juicers". The names are derived from two of the larger rental companies "Birds" and "Lime".



A successful hunt for this Bird-catcher

Both solutions get the job done, but in either case one may be faced with one of two problems:

- 1: A company is competing with other scooter-renting companies in several cities. In order to gain the edge, they want to improve the cost and time it takes to charge their scooters, such that they can be more quickly available for use.
- 2: A self-employed "Bird-catcher" will be competing with other freelancers on collecting these valuable "scooter-bounties". It is first come first served, and the faster they are collected, the better.

In both cases, a solution could be to plan a route for collecting such that the distance travelled is minimized. How should you go about planning a route such that you collect as many scooters as you can in the shortest amount of distance?

1.1 Problem introduction

The task laid out in front of us has an obvious appliance in the field of algorithms. It is quite similar to other wide-studied problems, one of which is the Traveling Salesman Problem. In the Traveling Salesman Problem (TSP) you are given a set of cities that you wish to visit, and the distance between each city. In what order should you visit these cities such that the total distance covered is minimized?

TSP is usually formulated as a graph problem. In the field of mathematics and algorithms a graph is a structure amounting to a set of objects in which pairs of objects are in some way related. An object is abstracted as a vertex, or node, and the relations between vertices are known as edges, or arcs. Graphs

can be weighted on the edges, where a number is assigned to each edge. In TSP, a city is notated as vertex, and the roads between the cities are the edges. The edges are weighted such that the value each edge is assigned represents the distance between two respective cities.

Our problem is similar, but differs in some areas; The cities are replaced by electric scooters, and you are not able to "visit" every scooter in one go. There could be hundreds of scooters scattered across the city, and only a selected amount can be collected before returning to charging station is necessary. The area of focus can be one of many: One is to find the shortest trip that visits 20 scooters and returns to base. The focus on this thesis, however, will be on calculating all trips such that every scooter is collected in one of the trips.

1.2 Problem preliminaries

In practical terms, we will not be referencing the scooter-collecting aspect much further, but move on to a more general approach, that can solve a variety of real-life problems, one of which could be scooter-collecting. The things in common for these problems are logistics-related:

- 1: They have a large number of Points-of-Interest (POI) with a known location, and a primary location (base), which will be our starting point.
- 2: The POI and base are spread across a plane and all POI should be both pair-wise reachable and reachable from base.
- 3: The POI are stationary, and some limitation is in place such that after a certain number of POI are visited, it is necessary to return to our primary base location.

In order to induce a graph problem from this, we need a data set with graph-like components, such as a road map. Where roads can be represented as edges, and intersections as vertices. POI located on this road-map should have coordinates, or a direct correlation to an intersection or road.

In our heuristic we can be given any graph with a strongly connected component¹ with at least 1 POI and 1 base. One vertex must be marked as the base. We expect both the edges and vertices to be weighted. The edge-weights will represent the distance between the end-points of the edge. And the vertex-weights represent the number of POI located on any given vertex. Depending on the type of graph, some pre-processing must be done. The base vertex is not considered in these steps:

- If an edge in the graph is undirected, we duplicate the edge and make them directed both ways.

¹A strongly connected component is a directed subgraph where all vertices are pair-wise reachable.

- If the graph is not connected, we induce our graph from the strongly connected component where the base is located.
- If the graph is not weighted on the edges, we assign each edge to have weight 1
- Vertices with weight more than 1 are split into new vertices with a new vertex for each weight value-point. The new vertices are connected with edges weighted 0 and share the same edges as the original vertex.
- If the graph is not complete², we calculate the shortest paths between all vertices with weight at least 1 and the base vertex. Then we induce a complete graph from these distances.
- If the graph is complete, but not without shortcuts³, we do the same as the previous point.

The reason for performing these steps is to streamline our algorithm such that a simple greedy heuristic can be abstracted. The new graph obtained after these steps can be traversed much simpler than the original graph:

- Step 1 is needed so that we do not need to differentiate between directed and undirected edges.
- Step 2 ensures that our graph is strongly connected.
- Step 3 ensures that each node needs only be visited once. It removes the unpleasant predicament of arriving at a vertex where the weight on the vertex is more than the remaining capacity for the current trip.
- Step 4 gives us a much simpler starting point as the structure of the graph is similar to the preferred graph type used in TSP solving techniques. The complete graph is simply a presentation of a matrix holding all the shortest paths between vertices.
- Step 5 is just an extra insurance so that we never need to think about finding the shortest path in the graph.
- Step 6 essentially "removes" the zero-weighted nodes by not including them in our complete graph.

After completing these pre-processing steps we end up with structurally speaking, homogeneous data types. Which is a directed strongly connected graph and a complete graph of all POI's located in the network.

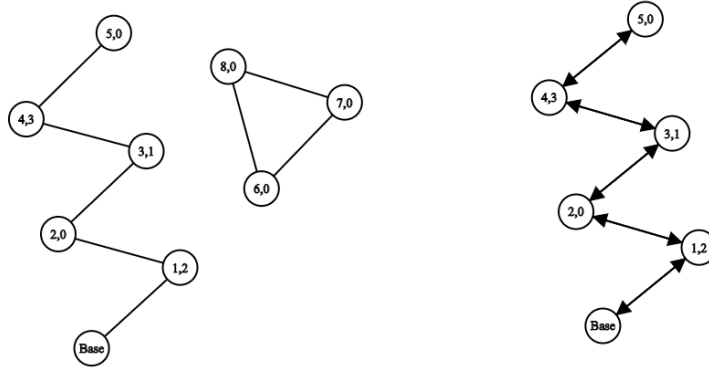
We expect the graph to be derived from a roadmap or a similar network, and in order to accomodate a varying format of the graph induced from the map data,

²A complete graph is a graph in which all vertices are pair-wise connected by an edge

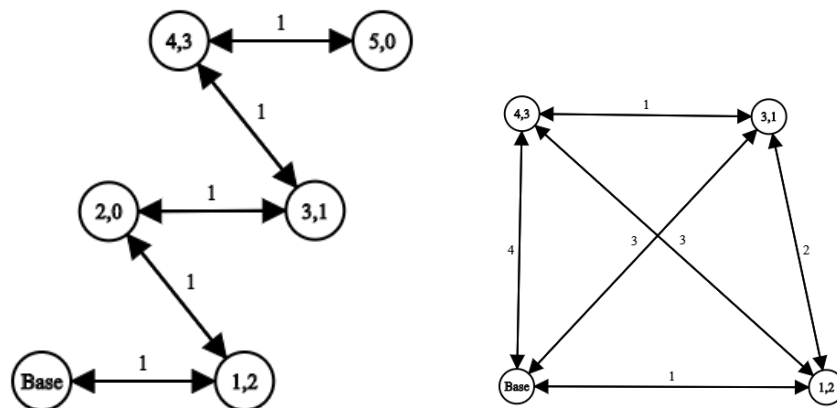
³A graph contains shortcuts if given 3 vertices $u, v, and w$ and a function $d(a, b)$ which indicates the distance between a and b : $d(u, v) + d(v, w) < d(u, w)$.

these steps are deemed as necessary. The data may be derived from a graph representing a road-network, which can be both undirected and directed (i.e. the existence of one-way streets), or it may only consist of vertices indicating the locations of our POI. Thus, we decrease the need to accomodate our algorithm to the input data.

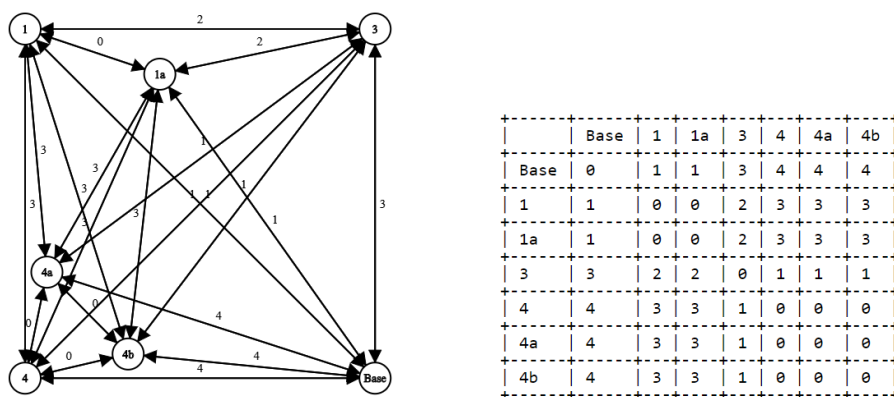
The order which these steps are performed in is somewhat arbitrary, but can affect the performance of the pre-processing heuristic. One might consider calculating the complete graph before splitting multi-weight vertices, as the distance data is identical for all duplicated vertices. The following example shows step-by-step actions taken when most cases apply:



Remove nodes not reachable from base, induce directed graph



assign edge-weights, Calculate complete graph,



duplicate heavier nodes, induce matrix for readability

2 Problem formulation

In this section we introduce two problems that are very similar but have one key difference. The first problem wishes to minimize both the length for visiting all nodes, and the number of trips needed to visit all.

Minimizing the number of trips is trivial; it can be calculated as the size of the graph divided by the number of nodes that can be visited in a single trip. The reason for demanding this will be explained further on.

The other problem presented is identical to the first, but there is no need to minimize the number of trips. This presents an interesting dilemma. In the first problem, as long as you can continue your trip, you must. But in the second, you must consider if you should terminate your trip, or continue. Developing a heuristic for the second problem will therefore be more difficult, but will likely yield a smaller minimization than the first problem.

2.1 Minimize sum of length of trips, Minimize number of trips

Given a complete weighted graph G , a start node s , a trip-length k ; Obtain a set of trips T which visits all of $V(G)$ such that each trip (except the last trip) visits exactly k nodes (excluding s) and the total length of all trips is minimized. A trip is a subset of vertices of G , which starts and ends with s . The length of a trip is the sum of all weights of edges between all adjacent vertices in the trip subset.

Input : A weighted graph G , vertex s , trip – length k

Problem : Minimize the length of the set of trips

$T : t_1, t_2, \dots, t_n$ where a trip $t_i \subseteq V(G)$

is an ordering of vertices v_1, v_2, \dots, v_m

s.t. the following holds for any trip t_i :

1. $\forall v \in V(G), v \in t_i \wedge v \notin T \setminus t_i$
 2. $\forall t_i \in T \ v_1 = v_m = s$
 3. $m = k + 2, |t_n| \leq k + |v_i = s|$
- (1)

The length of the set of trips is

$$\sum_{i=1}^n \sum_{p=1, q=2}^m w(v_p, v_q) \in t_i$$

This problem, we will call the Exact- k -Round Restricted Traveling Salesman Problem or EkRRTSP. A trip should minimize the total distance it travels. It is trivial to observe that this is an NP-Hard formulation as a specialized instance of this problem is the Traveling Salesman Problem. The problem becomes TSP if we set $k = |V(G)|$ and let the weight of all edges connected to s be weighted 0.

Note that we require each trip to visit the maximum number of nodes capable as long as there are at least k nodes left unvisited. This is a potential requirement to look at, as we may compare a heuristic that does not prohibit smaller trips and see if there is a possibility of getting a better performance.

2.2 Limit sum of length of trips, all trips are length at most k

Given a complete weighted graph G , a start vertex s , a trip-length k and bound X ; Is there a set of trips T which visits all of $V(G)$ such that the number of nodes visited excluding s is at most k and the total length of all trips is minimized?

Input : A weighted graph G , vertex s , trip – length k

Problem : Minimize the length of the set of trips

$T : t_1, t_2, \dots, t_n$ where a trip $t_i \subseteq V(G)$

is an ordering of vertices v_1, v_2, \dots, v_m

s.t. the following holds for any trip t_i :

1. $\forall v \in V(G), v \in t_i \wedge v \notin T \setminus t_i$
 2. $\forall t_i \in T, v_1 = v_m = s$
 3. $m \leq k + 2$
- (2)

The length of the set of trips is

$$\sum_{i=1}^n \sum_{p=1, q=2}^m w(v_p, v_q) \in t_i$$

This problem we call the Restricted Traveling Salesman Problem, or RTSP, and is essentially the same as the previous with the only difference being the length of each trip no longer requires a total visit of k nodes if possible. This is likely to never yield a worse solution than EkRRTSP as it is able to give the same solution as EkRRTSP.

3 Solving techniques

For solving the problems presented we will be using some common techniques that are further explained in this section.

3.1 Greedy algorithms

For EkRRTSP we will be using an algorithm that falls under the category of greedy algorithms. Greedy algorithms are characterized by one simple logical heuristic: At any given stage in our algorithm, select the path that is the local optimum.

This reduces the need to compute exhaustively on data, but can lead to locking out an optimal solution. As such, greedy algorithms have a varying degree of success on problems. In some problems, we find the global optimum with a

greedy heuristic. Others can approximate the optimum at a lower-bound guarantee. And some may yield the worst possible result.

A problem closely resembling EkRRTSP, TSP, has a greedy heuristic: Go to the currently closest unvisited city. This heuristic does not always find the global optimum for TSP, and in fact it is possible to obtain the worst possible solution using this technique. A similar technique shall be applied for EkRRTSP, as a baseline for other heuristics to compare to.

3.2 Clustering

In solving RTSP we will use a technique known as clustering. Clustering is commonly applied in the field of machine learning, as it is a tool for categorizing entities that have strong correlation. In our case the correlation is the distance between POI.

Clustering algorithms seek to group entities such that entities in the same group (or cluster) resemble each other more than entities of other groups. Cluster analysis may use a wide variety of measuring the likeness of entities, often more than a single measure.

For visualization purposes, the entities are placed on a plane where the closer two entities are to each other, the more they are alike. The likeness can therefore be measured as the distance between the entities. In RTSP the only measure of likeness will be the distance itself between them. So it lends an obvious appliance to the field of cluster analysis. Instead of categorizing entities in an unsupervised machine learning environment we use the technique to assign scooters to each trip.

There are many different types of clustering algorithms, which have their pros and cons. An important aspect that we need is the ability to limit the size of each cluster. As we want a cluster to be completed in a single trip. Hierarchical clustering techniques allow this.

3.2.1 Hierarchical clustering

Hierarchical clustering (HC) is a technique commonly used in machine learning. It is useful due to the fact that it does not need a specified number of clusters, as opposed to the more commonly used k-means clustering.

HC can be implemented in a top-down approach, where the initial cluster is all elements, and is continuously split into smaller clusters in a tree-like fashion until we are left with a single element in each cluster.

Or it can be implemented bottom-up, which is the complete opposite. We start with a single element cluster and then continuously merge clusters until we are left with a single cluster.

Both methods create a tree of clusters that applies itself to data of hierarchical properties. This is useful when categorizing entities that can be placed into several subcategories, such as animals (split/merge for each trait) or similar data elements. Although there are usually no categories for our POI, it is still possible to apply HC to our POI data.

In order to utilize HC for our problem, we will be starting with clusters containing a single POI and the base. Then, we merge clusters that have the lowest cost function (the closest ones) as long as the size of the merged cluster does not exceed our given limit value for size of each cluster. This will usually not yield a complete tree as usually the limit value is lower than the total number of POI.

Instead, we are left with several subtrees where each root is a cluster of size less than the limit. The idea is to then compute a TSP result on each root, and then moving onto their children. While the children give a better result combined than the parent, we keep moving down the tree. This will ultimately yield us a solution for solving RTSP using HC.

3.3 Dynamic Programming - rough draft

Once we have our clusters there still remains to calculate the shortest distance to visit each POI in the cluster. There is no way around accepting that this is the classic Traveling Salesman Problem mentioned several times in this thesis. The only difference is the base which we must start and stop our trip at. TSP is a notoriously difficult to solve, but there exists some techniques that outperform a brute force computation.

But let's first discuss how to consider our base node such that we guarantee that it is the first and last node in our trip. A simple way to do this is to duplicate the base and add a dummy node to each trip that have 0 weight on edges connected to its base and infinite weight to all other nodes. This will force the algorithm to start and end on the dummy nodes and thus, the base node. This may be a risky tactic as many TSP solvers do not handle the occurrence of 0 and infinite edges well.

4 Solution

First, in order to obtain the input given to our heuristic, we must perform some pre-processing steps that were explained in section 1.2

4.1 Pre-processing

The following pseudo-code tries to illustrate the pre-processing algorithm in a short concise matter. This should not be considered a direct representation of the logic in the source code.

4.1.1 Induce directed graph

First we simply want to be working with a directed graph. In the case of undirected edges in the graph we convert them into two opposite directed edges. After this procedure we consider all edges in the graph as directed.

Algorithm 1 Construct graph as directed

```
1: procedure DIRECTED_GRAPH(Graph  $G$ , node  $b$ )
2:   for each edge  $(u, v)$  in  $E(G)$  do
3:     if  $undirected(u, v)$  then
4:        $G.add\_edge(v, u)$ 
5:     end if
6:   end for
7:   return  $G$ 
8: end procedure
```

Runtime analysis:

We assume $undirected(u, v)$ is an $O(1)$ which yields us a runtime of $O(n)$ for n nodes.

4.1.2 Strongly Connected Component of Base node

Next, we consider graphs that are not strongly connected, and induce the sub-graph that is a SCC containing our base node. This is because we only want vertices that are strongly connected to the base vertex. To determine the SCC of b , we select the nodes which are reachable from b , and b is reachable from each respective node.

Algorithm 2 Remove vertices not strongly connected to node b

```

1: procedure SCC OF B(Directed graph  $G$ , node  $b$ )
2:    $reachable \leftarrow DFS(G, b)$   $\triangleright$  Run DFS to get all nodes reachable from  $b$ 
3:   for each node  $u$  in  $V(G)$  do
4:     if  $u$  in  $reachable$  then
5:        $connected \leftarrow DFS(G, u)$ 
6:       if  $b$  in  $connected$  then
7:         continue
8:       else
9:          $G \leftarrow G.remove(u)$ 
10:      end if
11:    else
12:       $G \leftarrow G.remove(u)$ 
13:    end if
14:  end for
15:  return  $G$ 
16: end procedure

```

Runtime analysis:

DFS is $O(n + m)$ and we potentially run DFS for each node, which yields a runtime of $O(n(n + m))$

4.1.3 Linkage matrix construction

The bulk of our pre-processing is to construct a linkage matrix for our clustering technique. We only consider nodes that are weighted, i.e. nodes that contain one or more POI's. This is also the computationally slowest part of the pre-processing as we need to make use of Dijkstra's algorithm on each weighted node to construct the matrix. It is possible to construct a linkage matrix that is sorted for each row. This is due to the nature of Dijkstra visiting nodes that are closest first. This will be helpful for our greedy heuristic. Our source code includes a non-sorting method as well, to maintain a structure of readability, i.e. the 4th element in a row always corresponds to the same POI.

Algorithm 3 Create linkage matrix for nodes that are weighted and b

```

1: procedure LINKAGE MATRIX(Directed graph  $G$ , node  $b$ )
2:    $M \leftarrow Matrix[V(G)][V(G)]$ 
3:    $dists \leftarrow djikstra(G, b)$  ▷  $dists[u]$ : distance from  $b$  to  $u$ 
4:   for each node  $u$  in  $V(G)$  do
5:     if  $w(u) > 0$  then
6:        $M[b][u] \leftarrow dists[u]$ 
7:        $udists \leftarrow djikstra(G, u)$ 
8:       for each node  $v$  in  $V(G)$  do
9:         if  $w(v) > 0$  or  $v == b$  then
10:           $M[u][v] \leftarrow udists[v]$ 
11:        end if
12:      end for
13:    end if
14:  end for
15:  return  $M$ 
16: end procedure

```

Runtime analysis:

The same principle as SCC applies, except we compute Dijkstra on each weighted node. Dijkstra's runtime complexity depends on the implementation, but optimally it will have a runtime of $O(n \log n)$. Which ultimately yields a runtime of $O(wn \log n)$ when constructing the linkage matrix.

4.1.4 Create a matrix row for each individual POI on each node

Finally, to simplify the clustering step, we seek to create individual rows for each POI. This means we need p duplicates for a row corresponding to a node weighted p . We don't need the duplicates in the graph, but they should be included in the linkage matrix.

Algorithm 4 Duplicate nodes equal to their weight, update M

```

1: procedure DUPLICATE_NODES(Graph G, linkage matrix M, node b)
2:   for each node  $u$  in  $V(G)$  do
3:     if  $w(u) > 1$  then
4:       for  $(i \leftarrow 1; i < w(u); i \leftarrow i + 1)$  do
5:          $M[ui] \leftarrow M[u]$   $\triangleright$  Update linkage matrix with new node
6:       for each row  $v$  in  $M$  do
7:          $v[ui] \leftarrow v[u]$ 
8:       end for
9:     end for
10:     $w(u) \leftarrow 1$ 
11:  end if
12: end for
13:  return  $M$ 
14: end procedure

```

Runtime analysis:

Although the concept seems rather simple, the current heuristic contains three nested for-loops. This is due to the fact that we must update the entire linkage matrix for each new row created. This procedure adds another $O(n^2W)$ to the total runtime for our pre-processing heuristic.

4.1.5 Runtime analysis of pre-processing

The challenge of the pre-processing is how to deal with the duplication of rows. The current heuristic will have a runtime of $O(n^2 \log n + n^2W)$. The total weight of the nodes play a role both in the pre-processing part as well as the later clustering procedure. It is currently unknown if there are methods to decrease the runtime of the pre-processing.

4.2 Greedy heuristic for EkRRTSP

ekRRTSP has a simple greedy heuristic: Simply visit the closest non-visited node at each step, and then return to base once k nodes are visited. The distance matrix is sorted by distance to easily find the closest vertex. The pseudocode below only returns both the total distance travelled to solve and the order of which the POI's are visited.

Algorithm 5 Greedy heuristic for EkRTSP

```
1: procedure GREEDY SOLVER(Sorted Linkage Matrix  $M$ , Size  $k$ , Start Node  
    $b$ )  
2:    $totalDist \leftarrow 0$   
3:    $trips \leftarrow [b]$  ▷ Order of nodes visited when solving  
4:    $visited[M[0]] \leftarrow all(False)$   
5:    $visited[b] := True$   
6:    $cnt \leftarrow 0$   
7:   while  $cnt < visited.size()$  do  
8:      $tripLength \leftarrow 0$   
9:      $currentNode \leftarrow b$   
10:    while  $trip < k$  do  
11:      if  $cnt \geq visited.size()$  then  
12:        break  
13:      end if  
14:      for ( $i \leftarrow 0$ ;  $i < M[currentNode].size()$ ;  $i \leftarrow i + 1$ ) do  
15:        if not  $visited[i]$  then  
16:           $visited[i] \leftarrow True$   
17:           $cnt \leftarrow cnt + 1$   
18:           $next \leftarrow i$   
19:          break  
20:        end if  
21:      end for  
22:       $trips \leftarrow trips.push(next)$   
23:       $totalDist \leftarrow totalDist + M[currentNode][next]$   
24:       $currentNode \leftarrow next$   
25:       $tripLength \leftarrow tripLength + 1$   
26:    end while  
27:     $trips \leftarrow trips.push(b)$   
28:  end while  
29:   $trips \leftarrow trips.push(b)$   
30:   $totalDist \leftarrow totalDist + M[currentNode][s]$   
31:  return ( $totalDist, trips$ )  
32: end procedure
```

4.2.1 Runtime analysis of greedy heuristic

While there are two nested while loops as well as a for loop the runtime is actually $O(n \log n)$. As each iteration of the inner for loop marks exactly one POI as visited. And the linkage matrix is sorted, so we need only find the closest node that is not visited. This is a very fast way to solve a TSP-like problem, but sadly not an effective one.

4.3 Agglomerative clustering heuristic for RBTSP

The heuristic for solving RBTSP using hierarchical clustering is split into three parts, where the two first parts can be considered pre-processing. The pseudocode snippets describe a slightly modified agglomerative clustering heuristic (ACH). We continuously select the cluster-pair that have the lowest cost function value and merge them. The cost-function is most commonly based on the distance-metric between clusters, we will be utilizing three of these in our experimentation:

- single linkage: The smallest distance between any two elements from each cluster. For our clusters we need to exclude the base vertex to be selected, as it is present in both clusters. This means that the absolute smallest distance between elements would be 0 (base element in both clusters). Single linkage tends to produce "loose" clusters.
- complete linkage: The largest distance between any two elements from each cluster. Complete linkage tends to create more compact clusters, and is the most common of the three.
- average linkage: The average distance between all pair-wise elements from each cluster.

Overall, complete linkage is the most likely to produce better results, but we include the other linkage methods for comparison's sake.

A regular ACH will keep merging clusters until a single cluster is obtained. However, in our heuristic we do not merge clusters if their combined size exceeds the imposed limit on cluster size. Thus, we are left with several trees, instead of a single one.

4.3.1 Initializing the clusterset

Our starting input is the linkage matrix obtained in the pre-processing section. We create a unique cluster for each row entry in the linkage matrix, which contains the respective vertex and the base vertex from the input graph. The base vertex is part of the linkage matrix, but does not get its own cluster.

Algorithm 6 Construct initial clusters

```
1: procedure INITIAL CLUSTERS(Linkage Matrix  $M$ , Start Node  $b$ )
2:    $clusterArray \leftarrow Array < cluster >$   $\triangleright$  Create initial clusters containing
   a single node and  $b$ 
3:   for  $e$  in  $M$  do
4:     if  $e \neq b$  then
5:        $elemList \leftarrow [b, e]$ 
6:        $cluster \leftarrow cluster(elemList)$ 
7:        $clusterArray \leftarrow clusterArray.push(cluster)$ 
8:     end if
9:   end for
10:  return  $clusterArray$ 
11: end procedure
```

Runtime analysis:

This procedure is completed using a single for loop over the rows of M , and so it is an $O(n)$ operation.

4.3.2 Initialize priority queue

Agglomerative clustering merges the two clusters with the lowest cost function. A fast way of finding these two clusters is to use a priority queue, or a heapqueue. We store the two clusters in the heap and sort them by the computed cost function belonging to the pair. We do not implement our own priority queue framework, instead using the library function included with Python 3.

Algorithm 7 Construct Priority Queue for pair-wise cluster distances

```
1: procedure PRIORITY QUEUE(list  $clusterArray$ , size  $k$ )
2:    $PQ \leftarrow priorityQueue()$ 
3:   for  $c1$  in  $clusterArray$  do
4:     for  $c2$  in  $clusterArray$  do
5:       if  $c1.size() + c2.size() \leq k$  then
6:          $res \leftarrow CostFunc(c1, c2)$ 
7:          $PQ \leftarrow PQ.push(res, (c1, c2))$ 
8:       end if
9:     end for
10:  end for
11:  return  $PQ$ 
12: end procedure
```

Runtime analysis:

Adding elements to a priority queue takes $O(\log n)$. As we need the pair-wise distances of all clusters we need to add elements n^2 times, yielding a run-time of $O(n^2 \log n)$

4.3.3 Agglomerative clustering

Finally, when the priority queue is ready, we can start building the trees. The clusters have a child and parent parameter that keeps track of the tree as it is being built. The parameters are updated in the *merge* function, where two clusters are merged into one. To make sure that each cluster is merged at most once, we include a lookup array that keeps track of which clusters have been merged. Once two clusters are merged, we need to recalculate the cost functions of the new cluster and all other non-merged clusters, and add them to the priority queue.

Algorithm 8 Pseudocode for Agglomerative clustering

```

1: procedure AGGLOMERATE(array clusterArray, priority queue PQ, size k)
2:   merged  $\leftarrow$  [false for c in clusterArray]
3:   while !PQ.isEmpty() do
4:     (dist, (c1, c2))  $\leftarrow$  PQ.pop()
5:     if !merged[c1] and !merged[c2] then
6:       newC  $\leftarrow$  merge(c1, c2)
7:       merged[newC]  $\leftarrow$  false
8:       merged[c1]  $\leftarrow$  true
9:       merged[c2]  $\leftarrow$  true
10:      for c in clusterArray do
11:        if !merged[c] then
12:          res  $\leftarrow$  costFunc(c, newC)
13:          PQ  $\leftarrow$  PQ.push(res, (c, newC))
14:        end if
15:      end for
16:      clusterArray  $\leftarrow$  clusterArray.push(newC)
17:    end if
18:  end while
19:  return clusterArray
20: end procedure

```

runtime analysis:

Priority queue is the aspect of importance when analyzing the runtime of the heuristic. Selecting and removing the smallest element of a priority queue is an $O(\log n)$ operation, as well as inserting new elements into the PQ. The initial size of PQ is n^2 so each of these operations become $O(\log^2 n)$. For each loop iteration, the number of clusters that remain unmerged is reduced by 1 (two clusters are merged, a new one is added to the pool). Also, the PQ does not shrink in size until the clusters are sufficiently large enough. Further analysis is required, but we know that the overall runtime is $O(n^2 \log n)$, for an $n \times n$ matrix which is slightly faster than not using a priority queue where the runtime would be $O(n^3)$.

4.3.4 Merge subroutine

To show that the set of clusters contain a tree-like structure we include a description of the *merge* sub-routine. Here we create the new cluster from two child-clusters and connect them with the child/parent parameters.

Algorithm 9 sub-routine for merging two clusters into one

```
1: procedure MERGE(cluster c1, cluster c2)
2:   elemList  $\leftarrow c1.elemList() \cup c2.elemList()$ 
3:   cluster newC  $\leftarrow cluster(elemList)$ 
4:   newC  $\leftarrow newC.setChildren((c1, c2))$ 
5:   c1  $\leftarrow c1.setParent(newC)$ 
6:   c2  $\leftarrow c2.setParent(newC)$ 
7:   merged(newC)  $\leftarrow$  false
8:   return newC
9: end procedure
```

runtime analysis:

The subroutine is a simple operation of connectors and object initiating, and has a runtime complexity of $O(1)$

4.4 Heuristic for solving TSP within each cluster, TBD

Finally we could use a traditional TSP solving technique for each group. Although the TSP solving step is exponentially hard to compute, so long as k is of a relative small size ($k < 30$) it should not be too resource intensive. As k is a constant it does not grow as the problem size grows so it is not considered an exponential solution for that matter.

Hopefully the clustering solution yields a smaller approximation factor than the proposed greedy solution while staying polynomial in time.

5 Results

5.1 Greedy heuristic

First we implemented a greedy heuristic with a simple logic. Once the pre-processing was complete, the solving was simple to implement. One obvious flaw with the heuristic that became apparent is that it does not take advantage of the split nodes. There is a high likelihood of having to visit the same split cluster several times to collect all the POI that are located there. This is something a more elegant solution may take better advantage of.

```

Input:
maximum trip: 2
B: B : 0.00, 0 : 1.00, 1 : 2.00, 2 : 3.00, 3 : 3.00, 4 : 2.00, 5 : 1.00,
0: B : 1.00, 0 : 0.00, 1 : 1.00, 2 : 2.00, 3 : 3.00, 4 : 2.00, 5 : 1.00,
1: B : 2.00, 0 : 1.00, 1 : 0.00, 2 : 1.00, 3 : 2.00, 4 : 3.00, 5 : 2.00,
2: B : 3.00, 0 : 2.00, 1 : 1.00, 2 : 0.00, 3 : 1.00, 4 : 2.00, 5 : 3.00,
3: B : 3.00, 0 : 3.00, 1 : 2.00, 2 : 1.00, 3 : 0.00, 4 : 1.00, 5 : 2.00,
4: B : 2.00, 0 : 2.00, 1 : 3.00, 2 : 2.00, 3 : 1.00, 4 : 0.00, 5 : 1.00,
5: B : 1.00, 0 : 1.00, 1 : 2.00, 2 : 3.00, 3 : 2.00, 4 : 1.00, 5 : 0.00,

Output:
total Distance: 15.0
trip 0 | length: 4.0 ['B', '0', '1', 'B']
trip 1 | length: 4.0 ['B', '5', '4', 'B']
trip 2 | length: 7.0 ['B', '2', '3', 'B']

Process finished with exit code 0

```

The image above is code input and output of the heuristic. 6 nodes excluding base. The input is the distance between nodes as well as the maximum trip possible.

The input table was derived from a graph that was originally an undirected cycle of 6 nodes including base, and a couple disconnected nodes that were removed in the pre-processing. The POI were generated at random, and it was possible for more POI to be assigned to the same node, leading to the node being split. However, in this example that was not the case, as there was a single POI for each node in the CC.

Output gives us the total distance covered and the trips performed, and the length of each trip. The total distance will be used as a benchmark to compare the different heuristics performance.