

Phase 1 Report

Matthew Ya & Aleece Randall

February 19, 2017

1 Quick How-to for Interface and Images

Blocked: Black
Regular unblocked: White
Hard to traverse: Orange
Regular unblocked with a highway: Cyan
Hard to traverse with a highway: Blue
Start: Lime Green
Path: Purple
End: Red

2 Implementation of Algorithms

We designed an abstract class called Search and 3 concrete classes UniformCostSearch, AStarSearch, and WeightAStarSearch which implemented the abstract class. Each concrete class shares the findPath() and findSuccessorSet() methods to traverse the map and expand the fringe. The abstract class Search has two abstract methods updateVertex() and setupFringe() that are implemented and defined by the concrete class. The reason for this is that each Search algorithm has its own way of calculating the value of $h(n)$, the value of $f(n)$ and setting up the fringe.

3 Optimization

The data structures used to implement the algorithms were chosen to achieve the best time and space complexities. The fringe or open list is a priority queue that is implemented using a binary heap. This data structure provides worst case complexities of $O(\log N)$ for search, inserts, and delete and $O(1)$ for removal of the cell with the lowest value. The successors of each node are stored in a HashSet because of its $O(1)$ insert. A HashSet was also chosen for the closed list after time/space complexity tradeoff. The HashSet has an average case search time of $O(1)$ but worst case time $O(n)$. The HashSet was chosen over a 2D boolean array due to the array's higher space requirement. Although, the Boolean array would yield a better worst case search time of $O(1)$, the space complexity is always the total number of nodes in the graph. On average, the HashSet is a much smaller value. The HashSet and Priority Queue were initialized to a higher but reasonable initial capacity in an attempt to reduce possible rehashing.

4 Heuristics

There are several possible ways to calculate the distance between the start node and goal node. Let dx represent the horizontal distance and let dy represent the vertical distance between the start node and goal node. Let D represent the vertical/horizontal cost and $D2$ represent diagonal cost. The costs used in the following heuristics focus on regular cells to try to avoid creating an inadmissible heuristic.

4.1 Euclidean Distance Formula

$$h = D2 * \sqrt{dx * dx + dy * dy} \quad (1)$$

Pro: considers diagonals

Con: computationally expensive, only considers a straight line

Using the diagonal cost, this function fails to take into account the faster speeds of the highways. Using an average of the highway horizontal/ vertical costs and unblocked cell diagonal s cost, the heuristic tries to consider the fact that an L-shaped movement on a highway is less costly than moving on a diagonal. We consider this to be the worst heuriistic.

4.2 Manhattan Distance Formula

$$h = D * (dx + dy) \quad (2)$$

Pro: computationally inexpensive

Con: does not consider diagonal movement

This formula considers the faster movements of highways. If the optimal path does not make use of the highways, this heuristic is not close to the actual cost. This heuristic is admissible in most cases, see below for exceptions.

4.3 Diagonal Distance Formula

$$h = D * (dx + dy) + (D2 - 2 * D) * \min(dx, dy) \quad (3)$$

Computes the number of steps if a diagonal is not taken and adds the minimum diagonal steps. The cost of the minimum diagonal steps is the net cost of not moving vertical/ horizontal. The equation tries to take into account both directions of travel. Considering all directions, highways for vertical/horizontal and regular cells for diagonal (extremely low probability for diagonal highway movement), might provide a good estimate that is admissible in most cases. However, using all the costs for a highway, this is the best optimal heuristic. It largely underestimates most paths but holds up for rare cases such where an optimal path is diagonal movements through parallel highways.

5 Experimental Results from 50 benchmarks for Phase 1

5.1 Heuristics

$$h1 = \sqrt{2} * \min(dx, dy) + \max(dx, dy) - \min(dx, dy) \quad (4)$$

$$h2 = 0.25 * (dx + dy) \quad (5)$$

$$h3 = ((0.25 + \sqrt{2})/2) * \sqrt{dx * dx + dy * dy} \quad (6)$$

$$h4 = 0.25 * (dx + dy) - (\sqrt{2} - 2 * 1) * \min(dx, dy) \quad (7)$$

$$h5 = \sqrt{2} * \sqrt{dx * dx + dy * dy} \quad (8)$$

$$h6 = .25 * (dx + dy) + (\sqrt{2} * .25 - 2 * .25) * \min(dx, dy); \quad (9)$$

5.2 Tables for Results

UniformCost Search	
Average Run Time	20.7ms
Average Path Length	180 nodes
Average Path Cost	295.54
Average Nodes Expanded	12636 nodes
Average Memory Used	5346.78 KB

A* Search						
Averages	h1	h2	h3	h4	h5	h6
Run Time	4.04ms	19.14ms	6.46ms	13.12ms	1.04ms	21.2ms
Path Length	153 nodes	180 nodes	162 nodes	171 nodes	134 nodes	180 nodes
Path Cost	354.71	295.54	306.97	304.01	400.76	295.54
Nodes Expanded	1202 nodes	8740 nodes	2957 nodes	5747 nodes	224 nodes	9600 nodes
Memory Used	431.78KB	3311.81KB	1101.76KB	2114.45KB	106.57KB	3871.26KB

Weighted A* Search with Weight 1.25						
Averages	h1	h2	h3	h4	h5	h6
Run Time	0.64ms	12.82ms	2.74ms	6.56ms	0.16ms	12.86ms
Path Length	132 nodes	178 nodes	154 nodes	164 nodes	126 nodes	178 nodes
Path Cost	402.93	295.98	342.58	321.03	439.24	295.89
Nodes Expanded	360 nodes	7518 nodes	1413 nodes	4388 nodes	137 nodes	8616 nodes
Memory Used	159.50KB	3017.65KB	431.34KB	1695.47KB	26.49KB	3504.37KB

Weighted A* Search with Weight 2.00						
Averages	h1	h2	h3	h4	h5	h6
Run Time	0.04ms	6.56ms	0.02ms	3.24ms	0.02ms	9.06ms
Path Length	125 nodes	168 nodes	126 nodes	141 nodes	123 nodes	174 nodes
Path Cost	463.33	301.86	436.88	400.35	484.69	298.39
Nodes Expanded	129 nodes	4541 nodes	149 nodes	2753 nodes	123 nodes	5904 nodes
Memory Used	53.06KB	1909.11KB	106.68KB	955.82KB	52.99KB	2442.98KB

6 Discussion

Overall, it's obvious that the most expensive search algorithm is the non-heuristic Uniform Cost Search algorithm. This algorithm has the highest average run time, nodes expanded, and memory used. The average performance of the A* search algorithm sits between the average performance of the Uniform Cost search algorithm and the average of the Weighted A* algorithm. For Weighted A* search algorithm, we noticed that the algorithm does significantly better in terms of memory and run time when the weight is set to 2 compared to when the weight is set to 1.25.

Out of the 5 tested heuristics, heuristic 5 had a positive influence in the behavior of the A* search algorithms. The heuristic with the worst result was heuristic 2. It seems that heuristic functions that come closest to predicting the actual path of two points will have a positive influence of the algorithm. Heuristic 5 is the least conservative heuristic function in regards to cost. Heuristic 5 provided the best computational performance but grossly overestimates the optimal cost.

In contrast, Heuristic 6 is an admissible/consistent function and a drastically worst performance. From these results, we can conclude that the more weight/value that originates from the heuristic leads to a better computational performance. The higher the heuristic values, the more influence it has in directing the algorithm. Heuristic 6 provides a fairly low estimate of the true optimal cost. The Weighted A* at 1.25 increased the computational performance without sacrificing the optimal cost nearly all of the trials. The Weighted A* at 2.00 more accuracy is lost. Heuristic 3 provides lower overestimates compared to Heuristic 5. As a result, the optimal path deviated a lot less than Heuristic 5 despite having the same Euclidean Distance Formula basis. The simulations demonstrate a tradeoff between computational performance and optimality.

7 Experimental Results from 50 Benchmarks for Phase 2

7.1 Heuristics

$$h1 = .25 * (dx + dy) + (\sqrt{2} * .25 - 2 * .25) * \min(dx, dy); \quad (10)$$

$$h2 = \sqrt{2} * \sqrt{dx * dx + dy * dy} \quad (11)$$

$$h3 = 0.25 * (dx + dy) - (\sqrt{2} - 2 * 1) * \min(dx, dy) \quad (12)$$

$$h4 = ((0.25 + \sqrt{2})/2) * \sqrt{dx * dx + dy * dy} \quad (13)$$

$$h5 = \sqrt{2} * \min(dx, dy) + \max(dx, dy) - \min(dx, dy) \quad (14)$$

7.2 Tables for Results

Sequential A* Search					
Averages	(1.25, 2.00)	(1.00, 2.00)	(1.50,2.00)	(1.25,1.50)	(1.25,2.50)
Run Time	53.64ms	74.04ms	45.62ms	92.84ms	25.78ms
Path Length	154 nodes	158 nodes	152 nodes	168 nodes	145 nodes
Path Cost	340.73	326.15	349.89	313.12	362.66
Nodes Expanded	6952 nodes	8871 nodes	5560 nodes	10810 nodes	4003 nodes
Memory Used	4002.58KB	4625.53KB	3002.39KB	4606.63KB	2450.08KB

Integrated A* Search					
Averages	(1.25, 2.00)	(1.00, 2.00)	(1.50,2.00)	(1.25,1.50)	(1.25,2.50)
Run Time	16.7ms	20.14ms	12.36ms	29.78ms	9.62ms
Path Length	163 nodes	169 nodes	156 nodes	171 nodes	154 nodes
Path Cost	326.75	315.87	337.61	307.78	348.71
Nodes Expanded	2076 nodes	2643 nodes	1804 nodes	4045 nodes	1177 nodes
Memory Used	1251.35KB	1630.06KB	1050.34KB	2135.95KB	847.06KB

8 Optimization

Similar to Phase 1, the closed lists in Phase 2 are implemented using a HashSet. The open lists were implemented using HashMaps with the nodes as keys and the algorithm-generated keys as values. At first, priority queues were used. When using priority queues, it updated the key value updates it across all the queues. It didnt change the positions in the queues but it created future issues with the comparator when the next value is inserted. Instead of cloning a new node for every queue, we switched the open lists to HashMaps. HashMaps provide an insert/update average time of $O(1)$. Unlike priority queues, the nodes do not need to be removed and re-inserted to update. The Top and MinKey methods have a time complexity of $O(n)$. For space efficiency, the Sequential Algorithm only initializes arrays (within the node object) that are needed to store the values for each heuristic for nodes that are expanded.

9 Proof for Sequential A* Search Question

Under the assumption that, $\text{Key}(s,0) \leq \text{Key}(u,0) \forall u \in \text{Open}_0$, then $g_0(s) \leq w_1 * c^*(s)$.

If we assume that $s = s_{goal}$, in other words we finally reach the goal node, then $g_0(s_{goal}) \leq w_1 * c^*(s_{goal})$. The anchor key $\text{Key}(s,0)$ represents the value that does not overestimate the cost from s to s_{goal} which means that the distance from s to s_{goal} is between $g_0(s_{goal})$ and $\text{Key}(s,0)$. This implies that the key of state s will not exceed the value $w_1 * c^*(s_{goal})$.

If Sequential Heuristic A* terminates in the i^{th} search, then we know that the following inequality was met (assuming $g_i(s_{goal}) < \infty$).

$$g_i(s_{goal}) \leq \text{OPEN}_i.\text{Minkey}() \leq w_2 * \text{OPEN}_0.\text{Minkey}().$$

We already know that $\text{OPEN}_0.\text{Minkey}()$ represents $w_1 * c^*(s_{goal})$ which bounds the anchor key of any state s . Which means we can rewrite the equation of $w_1 * w_2 * c^*(s_{goal})$ to $w_2 * \text{OPEN}_0.\text{Minkey}()$. This shows that the value $g_i(s_{goal})$ is bounded by $w_2 * \text{OPEN}_0.\text{Minkey}()$ which represents a sub-optimality factor of $w_1 * w_2$.

10 Proof for Integrated A* Search Question

The Integrated Algorithm guarantees that no state is expanded more than twice by maintaining two closed lists, an anchor list and a non-anchor list. A node is only expanded by a search if it is not a member of that closed list. If the anchor heuristic expands the node first, the node is never re-expanded because the algorithm checks if the node exists in the anchor closed list before checking the non-anchor list. The base condition to expand a node is that the g is lowered. Then if the node does not exist in the anchors closed list, it can be expanded by the anchor regardless if the node was already expanded in non-anchor list. Implementation of these properties can be found in the `expandState` method.