

AStar assignment — A routing problem

Lluís Alsedà

October 1, 2015

The assignment consists in computing an optimal path (according to distance) from *Basílica de Santa Maria del Mar* (Plaça de Santa Maria) in Barcelona to the *Giralda* (Calle Mateos Gago) in Sevilla by using the AStar algorithm. To this end one has to implement the AStar algorithm (compulsory in form of a function — not inside the main code) and compute and write the optimal path.

As the reference starting node for *Basílica de Santa Maria del Mar* (Plaça de Santa Maria) in Barcelona we will take the node with key (`@id`): 240949599 while the goal node close to *Giralda* (Calle Mateos Gago) in Sevilla will be the node with key (`@id`): 195977239.

For reference: My implementation (run in a very quick computer with a lot of memory) takes 20.51 CPU seconds to read and store the file and create the edges. The AStar computation itself takes 3.98 CPU seconds and the solution length is 959655.33 meters. It is shorter than the one found by Google since we are minimizing distances instead of time.

The map is contained in the file `spain.csv` that you have to download. This file has been created from `spain.osm` (written in *XML*) with the help of the *awk* program to get a format more friendly and easier to read than the original XML format.

The file `spain.csv` has the character `'|'` as field separator and it has three types of fields: `node`, `way` and `relation`. The structure of each field is:

```
node|@id|@name|@place|@highway|@route|@ref|@oneway|@maxspeed|node_lat|node_lon
way|@id|@name|@place|@highway|@route|@ref|@oneway|@maxspeed|membernode|...
relation|@id|@name|@place|@highway|@route|@ref|@oneway|@maxspeed|rel_type|type;@id;@role|...
```

Nodes specify a single point. *The list of nodes is sorted with respect to keys (@id)*. A way is a list of nodes (between 2 and 2000) and specify the relations between them (edges in the graph). Relations are mainly for drawing the maps and are irrelevant to our problem. Some data about the file:

- Longest line: 79857 chars
- Maximum number of fields: 5306
- Maximum width of `@name`: 184 chars
- Number of nodes: 23895681
- Number of ways: 1417363
- Number of relations: 25394533

The @id field has maximum length 10 chars. For comparison speed I recommend to store it as unsigned long.

The @oneway field takes two values: empty or oneway. If the pair of nodes A|B appears in the nodes list of the way, then in the graph always there is an edge from A to B. If the value of @oneway is empty (twoways) then additionally, in the graph there is an edge from B to A.

Warning: Unfortunately, the file is not consistent. There are ways with less than two nodes (that have to be discarded) and there are nodes in ways that do not appear in the list of nodes. They have to be omitted and the process of assigning edges to every pair of consecutive nodes in a way must be restarted.

Due to the concrete values of the @id's of nodes and the jumps between consecutive @id's (in the file), it is not feasible to use the @id's as indexes in the vectors. There is not enough memory (in the universe?) for that. Thus, I *strongly recommend to store the nodes in a vector of node structures and refer to the nodes internally in the program by the index in this vector* (not by the @id). This is way of naming nodes allows a quicker way of finding them. The price to pay for this is that at reading/storing time the arrows in the graph (specified in the way commands) have to be converted from @id numbering to vector index numbering.

I am using the following structure although everything (specially the adjacency relations can be implemented in other ways):

```
typedef struct {
    unsigned long id;           // Node identification
    char *name;
    double lat, lon;           // Node position
    unsigned short nsucc;      // Node successors: wighted edges
    unsigned long *successors;
} node;

typedef char Queue;
enum whichQueue {NONE, OPEN, CLOSED};

typedef struct {
    double g, h;
    unsigned long parent;
    Queue whq;
} AStarStatus;
```

The number of nodes is the dimension of this vector. To determine it for any data file it may be necessary to do a first reading of the file to compute this number (although it is given above for the file `spain.csv`).

An information that may help in deciding the storage model for the neighbours: the maximum valence in the graph is 16 (maximum number of nodes connected to a given one) and the average valence is 1.99. At reading time, to save `realloc`'s it might be useful to have an additional vector `unsigned short nsuccdim`; such that `nsuccdim[i]` contains actual dimension of `node[i].sucessors` (and it is initialized to zero).

Concerning the conversion of the arrows from @id numbering to vector index numbering: When processing the way's one has two keys A|B and has to search in the vector of nodes for the indexes (say a|b)

of the nodes with those keys to establish the edges between them (`node[a].sucessors[nsucc] = b`). This is a painful process where brute force does not work. For example, the loading of the map of Catalonia takes more than two hours by using this strategy. Thus, some efficient search method as binary search must be used. You may want to read from page 33 to page 47 (specially pages 39 and 41) of my notes (in Catalan) *Estructures i tipus de dades en C* that you can find in http://mat.uab.cat/~alseda/MatDoc/C_Estructures.pdf

An efficient approach to the problem (similar to what happens in GPS industry) is to write two programs:

- One that reads the file and computes the graph with binary search and stores the graphs in a *binary file* with arrows between nodes already determined. This file is thus very quick to read. To avoid a big number of `alloc`'s I recommend to store all names of all node in a single vector and all successors of all nodes in a single vector and store them in file like this.
- Then, a second program reads this formatted binary file, thus getting the map already in graph form, asks for the start and goal nodes and performs AStar algorithm. After reading the (big) vectors containing all names and successors, each node can set pointers to the parts of the vectors corresponding to this node. The reading of the binary file takes 0.76 CPU seconds in my implementation.

It is very important the choice of the heuristic distance (the concrete result strongly depends on this). Assuming that you choose as heuristic the shortest straight distance on the earth surface between two points, there are different ways of computing this distance (it is not well defined). Most common distances are Haversine formula (great-circle distance between two points), Spherical Law of cosines and others. I recommend you tyo consider carefully the distance you adopt. To this end, you may want to look at <http://www.movable-type.co.uk/scripts/latlong.html>