

Distance algorithms discription

Introduction

In Europe (see below for the delimitations chosen), there are over 65700 cities. This makes over $2.1 * 10^9$ combinations, which **is too much** to be calculated and filled into the database.

Furthermore, given the problem with Graphhopper described at the end of this document, using this directly would create errors for **nearly 45%** of combinations.

Base algorithm

Preparation

- First, you need a way to define **zones** : a delimited area with a given centroid to separate postcodes.
- Second, fill a database with all the distance between each pair of **zone**, by route and as the crow flies.

Get distance between postcodes

Given two postcodes (postcode, country iso code) :

1. Get the distance as the crow flies between the two postcodes (named C)
 2. Find to which group belong each postcode
 3. Extract from the database the two kinds of distance between the two *zones* and calculate the ratio $route/crow$ (named R).
 4. Finally calculate : $C * R$
-

2-digits version

The orinal proposal was to use the 2 first digits of the postcodes to define the *zones*. By "*2 first digits*", I think was meant : *Administrative region* (such as *Bundesland* in Germany, *Département* in France or *State* in the U.S.A.) ; and indeed in France and Germany, these *zones* are indicated by the two first digits of the postcodes.

Although these regions are often used for shipping, transport and administration, a few reasons could create inaccuracies with this approach :

- Not every country use the two first digits for their regions (for instance : if a country uses 3 digits and 305 is top north and 307 in the south, they would all be grouped under 30 in the middle of the country)
- Regions don't have the same size in every countries, this could lead in very different accuracy depending on the country
- Regions' shapes of the **administrative regions** aren't best suited for the distance of the centroids. (It is better if they are somewhat **round**).

I looked for a way to define administrative regions and their contained postcodes but couldn't find a proper solution for that. Finally, to implement this version of the algorithm, I used grouped every city with their forst

two digits and country code, **even though this doesn't necessarily make sense for all countries, and can cause inaccuracies.**

To get the centroid of these *zones*, I just calculated the average of the postcodes coordinates. This is maybe not a good way, but I don't see any other solution.

Using this method, we get **1840** *zones*, which makes **1 691 880** combinations.

Grid version

To solve the problems stated above (and make the development easier), I thought of using a **grid** instead of the 2-digit region. The grid is made by *meridians* and *circles of latitude*, at a certain step (every degree, tenth of degree, ...).

This solution has these advantages :

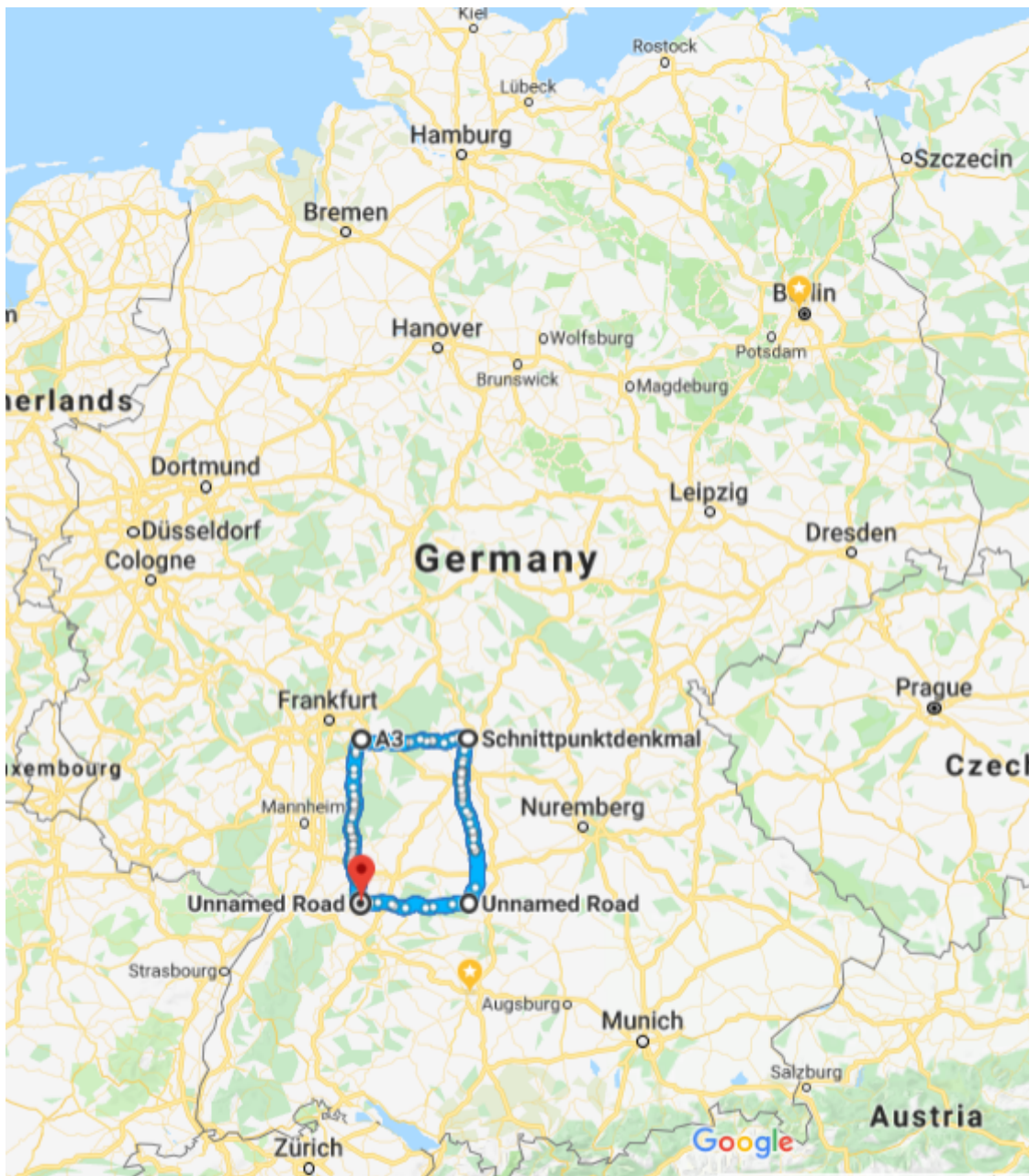
- It is much easier to define in which region is located each postcode (directly look at the coordinates)
- The *zones*' shapes are standardized and *approximatly* the same size.
- The *zones* have a *squary* shape which could provide better accuracy.

Coordinates in Europe (continental + Ireland and UK + Mediterranean Isles) roughly goes from 71°N,11°W (Norway and Ireland) to 35°N,41°E (Spain and Ukraine). So we have a difference of latitude and longitude : 36 and 52.

If we cut every degree, we have a total of **1872** *rectangles*. So we would have as a maximum number of distances between each couple of *rectangle* of $\binom{1872}{2} = 1,75 * 10^6$. By filtering out the *rectangles* that don't contain any postcodes, we have only **1056** rectangles, which gives **557040** combinations.

If we cut every tenth of degree, we have a total of **187200** *rectangles*, **29 010** only containing postcodes, which gives **420 775 545** combinations.

To get an idea of the size of such *rectangle*, here is an example in Germany between 49°N,9°E and 50°N,10°E, it is 85km wide and 130km long:



Graphhopper problem

Graphhopper has sometimes trouble tracing roads between two points when one of the coordinates are located not enough close from the road or something, **the request returns an error**.

This creates a problem when requesting the distance by road between two *zones'* centroids : an error occurs and the corresponding line is not entered in the database. For this reason, calculating the distance using any of the algorithms from the database won't work.

Solution

Probably one of the best solutions would be to *circle around* the centroids until we find coordinates that work, that would increase a lot the time spent generating the distance database and it would be difficult to implement.

For these reasons, when these errors happened, I first chose to default to the distance *as the crow flies*. After some tests, I found that, in average, the distance *as the crow flies* is **76%** of the distance by route. So I set the

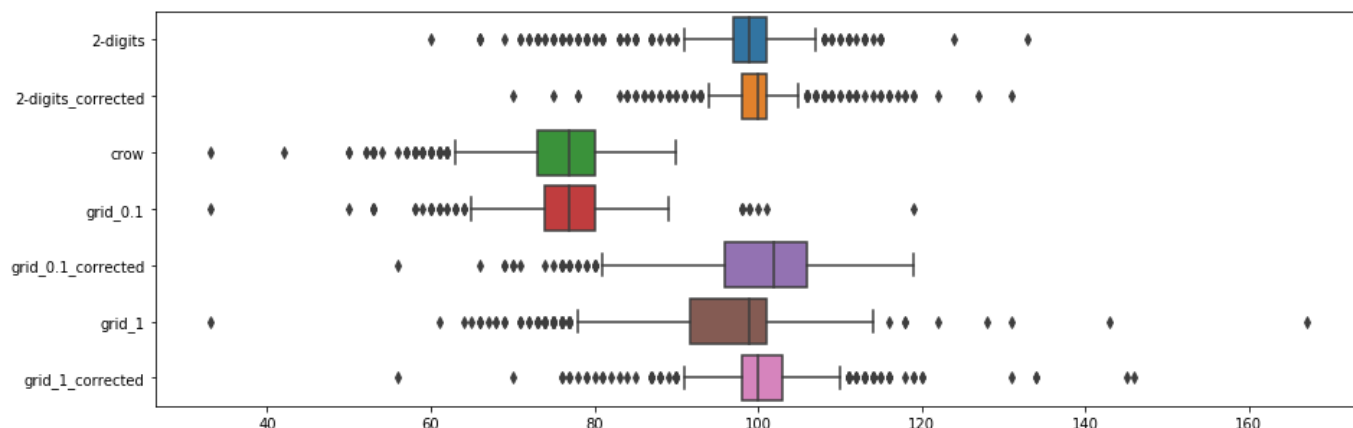
default to **dist_crow / 0.76**.

Note : This is also applied when the two requested postcodes are **in the same zone**.

Comparison on the accuracies of each methods

I ran the different distances algorithms on 1000 randomly picked pair of postcodes. On those, **Graphhopper returned an error for 45% of them**.

Here are the accuracies comparison on the 55% remaining, the values are in percentage of the route distance directly via Graphhopper :



Note : The Grid 0.1 database being not filled yet, its results are mostly the default to the crow and crow corrected.

- We can see that the 2-digits algorithm provides much better results than the grid 1-degree one. This can be due to the Graphhopper error and to the fact that many 2-digits regions in Europe are smaller than the 1-degree squares.
- The correction applied to solve the Graphhopper problem is quite efficient: all averages are around 100% and the data is concentrated on it.