

## Finding Gaps Between Mass Peaks - Analysis and Design Decisions

We have a list of mass spectrometry readings, with a mass value and an intensity value. First, we want to remove all low-intensity readings. Then, we want to find gaps in the masses that correspond to the masses in some table of compounds (with some tolerance, either percentile or absolute). Then we want to find all sequences of masses where one mass begins on a peak as another ends, with no subsequences. This is just a matter of comparing masses then checking if they lie within some threshold.

This problem can be modelled as a DAG (directed acyclic graph), where each vertex is enumerated and vertices can only have edges to vertices with higher numbers, edges have a label (and a given vertex may have several edges with different labels to another, although this is unlikely in practise), and these edges are unknown, but discoverable.

In this case, vertices are individual mass spectrometry readings, and edges are potential compounds from gaps between mass peaks.

There are multiple design decisions to be considered here. Firstly, the most naive algorithm treats any vertex as a potential beginning point for a sequence tag. This means for any vertex  $k$ , in the worst-case (which is very unlikely) we must check  $(k - 1) + (k - 2) + \dots + 2 + 1$  vertices. This is  $O(n^2)$ .

However if we visit a later vertex as part of discovering an earlier one, we know that any sequence tag that we visit as part of another one must be a subsequence. Then we can ignore these vertices as starting points. Then each new  $k$ th input can only be visited at most  $c(k - c - 1)$  times, where  $c$  is the number of previously visited nodes s.t.  $0 < c < k - 2$ .

...  
This resolves to  $O(n)$ .

Furthermore, we can search this graph either by first compiling all discoverable edges into some auxiliary data structure and then walking across it, or we can construct walks as we discover edges. While the auxiliary data structure method can be done by either breadth-first search or depth-first, it is slightly more naturally predisposed to breadth-first search, and the latter almost exclusively favours depth-first. The auxiliary data structure of course has the disadvantage of having to traverse it, so is slightly slower in best-cases, but prevents recalculation so should be significantly faster in worst-cases (and it also lets us maintain a more concise stack and fewer tables of mass differences). Traversing this auxiliary data-structure is also  $O(n)$ , by the same logic as above (it is essentially a model of the same graph with the edges discovered) but it should be significantly sparser than the data itself.

Thirdly, we can either discover edges by naive comparison with a Python for loop (saving a list of differences to prevent recalculation for each compound) or by mass-calculation with NumPy arrays. The latter should be significantly faster, but we shall benchmark this difference.