## 2.4   Compact storage of similar sequences

Consider the problem of storing a large set of strings, i.e., sequences of characters from a finite alphabet. We assume that the strings have many similar entries (they differ only in a small number of positions) and we wish to store them in a compact way. This problem arises in several contexts such as when storing DNA sequences, where the characters correspond to the four DNA bases. In this exercise, we consider the simplified version of the problem with only two characters.

Given a set of $k$ sequences of $M$ bits, we compute for each pair $i$, $j$, with $1 \leq i, j \leq k$, the Hamming distance between the sequences $i$ and $j$, i.e., the number of bits that need to be flipped in sequence $i$ to obtain sequence $j$. This function clearly satisfies the three usual properties of a distance: nonnegativity, symmetry and triangle inequality.

Consider the following set of 6 sequences and the corresponding matrix $D = \{d_{ij}\}$ of Hamming distances

<table>
<tr><td>1) 011100011101</td><td></td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td></tr>
<tr><td>2) 101101011001</td><td>1</td><td>0</td><td>4</td><td>4</td><td>5</td><td>4</td><td>3</td></tr>
<tr><td>3) 110100111001</td><td>2</td><td></td><td>0</td><td>4</td><td>3</td><td>4</td><td>5</td></tr>
<tr><td>4) 101001111101</td><td>3</td><td></td><td></td><td>0</td><td>5</td><td>2</td><td>5</td></tr>
<tr><td>5) 100100111101</td><td>4</td><td></td><td></td><td></td><td>0</td><td>3</td><td>6</td></tr>
<tr><td>6) 010101011100</td><td>5</td><td></td><td></td><td></td><td></td><td>0</td><td>5</td></tr>
<tr><td></td><td>6</td><td></td><td></td><td></td><td></td><td></td><td>0</td></tr>
</table>

where, due to symmetry, only the upper triangle of the matrix is shown.
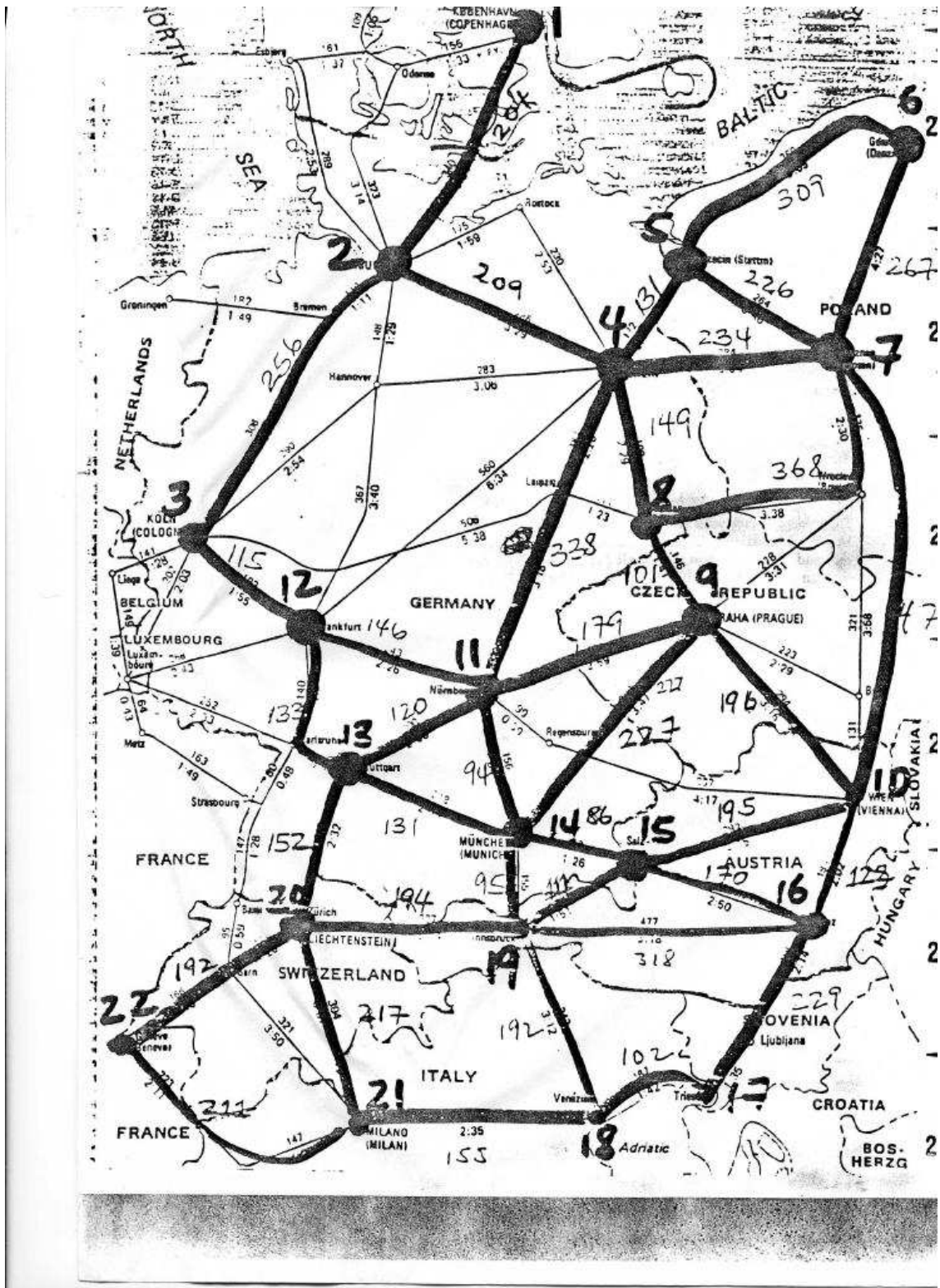
In order to exploit redundancies between sequences and to save memory, we can store: i) one of the sequences, called the reference sequence, completely and ii) for every other sequence, only the set of bit flips that allow us to retrive it either directly from the reference sequence or from another sequence.
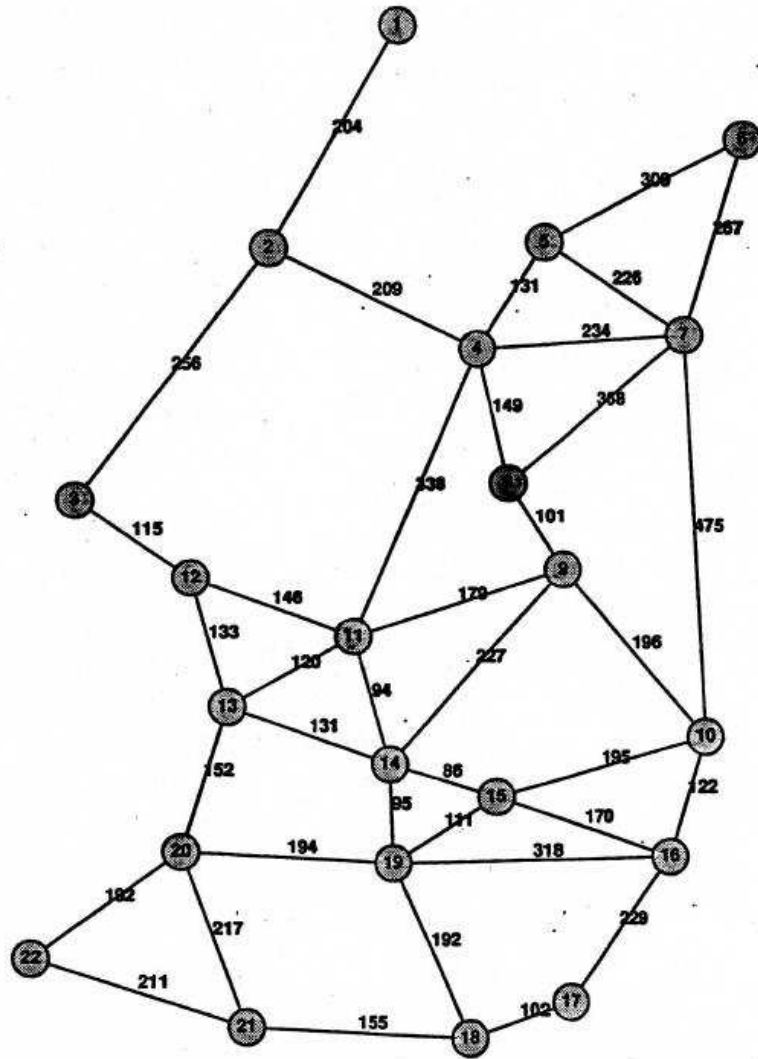
Show how the problem of choosing which differences to memorize, so as to minimize the total number of bits used for storage, can be reduced to the problem of finding a minimum-cost spanning tree in an appropriate graph. Solve the problem for the given instance.

[*Hint*: How many bits are needed to store the list of the differences between any given pair of sequences $i$ and $j$?]

## 2.5   Minimum cost network design

Given the following undirected graph representing the possible fiber optic connections between some major European cities together with the corresponding distances. Since the fiber optic connection cost is proportional to the length of the connection, design a network of minimum total cost which guarantees that every pair of cities is connected.
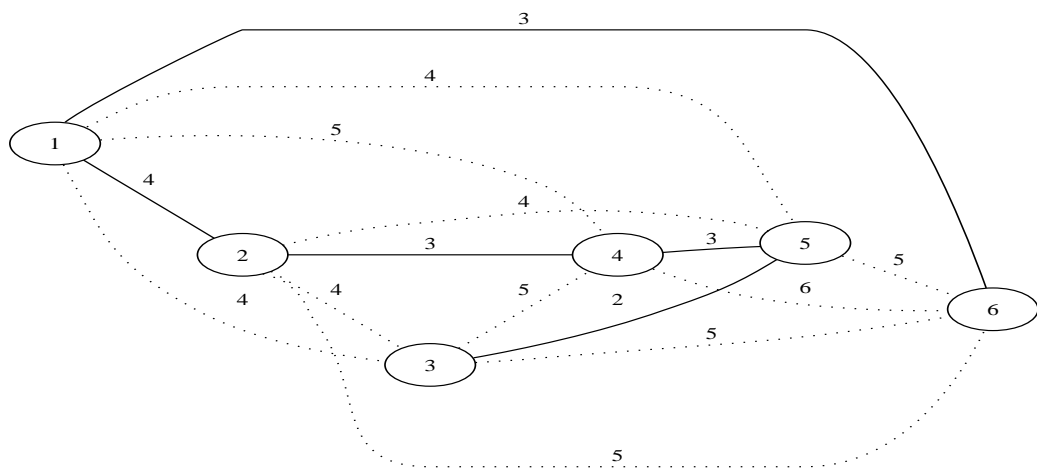
SOLUTION

2.4 **Compact storage of similar sequences**. We construct a complete graph $G$ with a node for each sequence and an edge for each pair of sequences. Moreover, each edge $(i, j)$ is assigned a cost $d_{ij}$.

The problem is then to look for a subgraph $G'$ of $G$ of minimum total cost. Since only one sequence is completely stored, $G'$ must be connected to be able to construct any sequence. Since a subgraph of minimal cost is sought, $G'$ will be acyclic. It followsa minimum-cost spanning tree in $G$ provides an optimal solution to the problem under consideration.

The following minimum cost spanning tree has been found using Prim's algorithm.



2.5 **Minimum cost network design**. Apply either Prim's or Kruskal's algorithm to the given instance.