# Worksheet 04

Name: Jeong Yong Yang, Junho Son UID: U95912941, U64222022

## Topics

- Distance & Similarity

## Distance & Similarity

Part 1

a) In the minkowski distance, describe what the parameters p and d are.

The parameter p is tune distance (parameter that we are able to tune). If p = 2, it is Euclidean distance and if p = 1, it is Manhattan distance.

The parameter d is the dimension. In other words, it represents d-dimensional space.

b) In your own words describe the difference between the Euclidean distance and the Manhattan distance.

Manhattan distance measures the distance horizontally and vertically but cannot go diagonal so when measuring distance between two points, it has to add the horizonal and vertical aspects. However, Euclidean distance can measure the distance diagonally and can measure the direct line distance from two points. Therefore, Manhattan distance is always greater than or equal than the Euclidean distance.

Consider A = (0, 0) and B = (1, 1). When:

- p = 1, d(A, B) = 2
- p = 2, d(A, B) = $\sqrt{2} = 1.41$
- p = 3, d(A, B) = $2^{1/3} = 1.26$
- p = 4, d(A, B) = $2^{1/4} = 1.19$

c) Describe what you think distance would look like when p is very large.

As p goes very large, the distance between A and B decreases and would be closer to 1 since it becomes 2^(very small number). In other words, as p increases, the distance between A and B goes closely to 1 but never reaches 1.

d) Is the minkowski distance still a distance function when p < 1? Expain why / why not.

No it cannot be a distance function when p < 1 because once p < 1, it violates the triangle inequality rule of d(i, j) <= d(i, k) + d(k, j).

e) when would you use cosine similarity over the euclidan distance?

We would use cosine dissimilarity over the Euclidean distance when direction matters more than magnitude.

f) what does the jaccard distance account for that the manhattan distance doesn't?

Jaccard distance accounts for the size of intersection, which Manhattan distance does not take into account.

## Part 2

Consider the following two sentences:

```
1 s1 = "hello my name is Alice"
2 s2 = "hello my name is Bob"
```

using the union of words from both sentences, we can represent each sentence as a vector. Each element of the vector represents the presence or absence of the word at that index.

In this example, the union of words is ("hello", "my", "name", "is", "Alice", "Bob") so we can represent the above sentences as such:

```
1 v1 = [1,     1, 1,     1, 1,     0]
2 #      hello my name is Alice
3 v2 = [1,     1, 1,     1, 0, 1]
4 #      hello my name is     Bob
```

Programmatically, we can do the following:

```
1 corpus = [s1, s2]
2 all_words = list(set([item for x in corpus for item in x.split()]))
3 print(all_words)
4 v1 = [1 if x in s1 else 0 for x in all_words]
5 print(v1)

    ['my', 'name', 'is', 'Bob', 'Alice', 'hello']
    [1, 1, 1, 0, 1, 1]
```

Let's add a new sentence to our corpus:

```
1 s3 = "hi my name is Claude"
2 corpus.append(s3)
```

a) What is the new union of words used to represent s1, s2, and s3?

```
1 union_of_words = list(set([item for x in corpus for item in x.split()]))
2 print(union_of_words)

    ['my', 'name', 'Claude', 'hi', 'is', 'Bob', 'Alice', 'hello']
```

b) Represent s1, s2, and s3 as vectors as above, using this new set of words.

```
1 vector1 = [1 if x in s1 else 0 for x in union_of_words]
2 vector2 = [1 if x in s2 else 0 for x in union_of_words]
3 vector3 = [1 if x in s3 else 0 for x in union_of_words]
4
5 print(vector1)
6 print(vector2)
7 print(vector3)

    [1, 1, 0, 0, 1, 0, 1, 1]
    [1, 1, 0, 0, 1, 1, 0, 1]
    [1, 1, 1, 1, 1, 0, 0, 0]
```

c) Write a function that computes the manhattan distance between two vectors. Which pair of vectors are the most similar under that distance function?

```
1 def manhattanDistance(vec1, vec2):
2   value = 0
3   for x in (range(len(vec1))):
4     value += abs(vec1[x] - vec2[x])
5   return value
6
7 man1 = manhattanDistance(vector1, vector2)
8 man2 = manhattanDistance(vector1, vector3)
9 man3 = manhattanDistance(vector2, vector3)
10
11 print(man1)
12 print(man2)
13 print(man3)

    2
    4
    4
```

Vector1 and vector2 (vector from "hello my name is Alice" and vector from "hello my name is Bob", respectively) are most similar under Manhattan distance.

d) Create a matrix of all these vectors (row major) and add the following sentences in vector form:

- "hi Alice"
- "hello Claude"
- "Bob my name is Claude"
- "hi Claude my name is Alice"

- "hello Bob"

```
1 sen1 = "hi Alice"
2 sen2 = "hello Claude"
3 sen3 = "Bob my name is Claude"
4 sen4 = "hi Claude my name is Alice"
5 sen5 = "hello Bob"
6
7 unique_words = list(set([item for x in corpus for item in x.split()]))
8 print(len(unique_words))
9 print(unique_words)
```

```
8
['my', 'name', 'Claude', 'hi', 'is', 'Bob', 'Alice', 'hello']
```

```
1 #      my,    name,   Claude,   hi,    is,    Bob,    Alice,   hello
2 vec1 = [0,    0,      0,        1,     0,     0,      1,       0]
3
4 vec2 = [0,    0,      1,        0,     0,     0,      0,       1]
5
6 vec3 = [1,    1,      1,        0,     1,     1,      0,       0]
7
8 vec4 = [1,    1,      1,        1,     1,     0,      1,       0]
9
10 vec5 = [0,   0,      0,        0,     0,     1,      0,       1]
11
12 matrix_of_vecs = [vec1, vec2, vec3, vec4, vec5]
13 print(matrix_of_vecs)
```

```
[[0, 0, 0, 1, 0, 0, 1, 0], [0, 0, 1, 0, 0, 0, 0, 1], [1, 1, 1, 0, 1, 1, 0, 0], [1, 1, 1, 1, 1, 0, 1, 0], [0, 0, 0, 0, 0, 1, 0, 1]]
```

e) How many rows and columns does this matrix have?

```
1 print(f"There are {len(matrix_of_vecs)} rows in the matrix.")
2 print(f"There are {len(matrix_of_vecs[0])} columns in the matrix.")
```

```
There are 5 rows in the matrix.
There are 8 columns in the matrix.
```

f) When using the Manhattan distance, which two sentences are the most similar?

```
1 sim_vec1 = []
2 sim_vec2 = []
3 index1 = 0
4 index2 = 0
5 smallestManD = manhattanDistance(vec1, vec2)
6 for vv1 in range(len(matrix_of_vecs)):
7   for vv2 in range(vv1+1, len(matrix_of_vecs)):
8     if manhattanDistance(matrix_of_vecs[vv1], matrix_of_vecs[vv2]) < smallestManD:
9       index1 = vv1 + 1 # add 1 since index starts from 0
10      index2 = vv2 + 1 # add 1 since index starts from 0
11      sim_vec1 = matrix_of_vecs[vv1]
12      sim_vec2 = matrix_of_vecs[vv2]
13      smallestManD = manhattanDistance(matrix_of_vecs[vv1], matrix_of_vecs[vv2])
14 print(sim_vec1, sim_vec2)
15 print(index1, index2)
```

```
[0, 0, 1, 0, 0, 0, 0, 1] [0, 0, 0, 0, 0, 1, 0, 1]
2 5
```

Vector2 and vector5 (vector from "hello Claude" and vector from "hello Bob", respectively) are most similar under Manhattan distance.

## ⌄ Part 3 Challenge

Given a set of graphs $\mathcal{G}$, each graph $G \in \mathcal{G}$ is defined over the same set of nodes $V$. The graphs are represented by their adjacency matrices, which are 2D arrays where each element indicates whether a pair of nodes is connected by an edge.

Your task is to compute the pairwise distances between these graphs based on a specific distance metric. The distance $d(G, G')$ between two graphs $G = (V, E)$ and $G' = (V, E')$ is defined as the sum of the number of edges in $G$ but not in $G'$, and the number of edges in $G'$ but not in $G$. Mathematically, this can be expressed as:

$$d(G, G') = |E \setminus E'| + |E' \setminus E|.$$

## ⌄ Requirements:

1. **Input**: Should take a list of 2D numpy arrays as input. Each array represents the adjacency matrix of a graph.

2. **Output**: Should output a pairwise distance matrix. If there are $n$ graphs in the input list, the output should be an $n \times n$ matrix where the entry at position $(i, j)$ represents the distance between the $i^{th}$ and $j^{th}$ graph.

```
1   def graph(graph):
2      num_graphs = len(graph)
3      distance_matrix = [[0 for _ in range(num_graphs)] for _ in range(num_graphs)]
4      for x in range(len(graph)):
5        for y in range(x+1, len(graph)):
6          distance = 0
7          for r in range(len(graph[x])):
8            for c in range(len(graph[x][r])):
9              if graph[x][r][c] != graph[y][r][c]:
10                distance += 1
11          distance_matrix[x][y] = distance
12          distance_matrix[y][x] = distance
13      return distance_matrix
14
15   gra1 = [[1, 1, 1],[0, 0, 1], [1, 0, 1]]
16   gra2 = [[1, 0, 0],[1, 1, 0], [0, 1, 1]]
17   gra3 = [[0, 0, 1],[0, 1, 0], [1, 1, 0]]
18
19   graphs = [gra1, gra2, gra3]
20   print(graph(graphs))
21
22   gra1 = [[1, 1, 1],[0, 0, 1], [1, 0, 1], [0, 0, 1]]
23   gra2 = [[1, 0, 0],[1, 1, 0], [0, 1, 1], [1, 1, 1]]
24   gra3 = [[0, 0, 1],[0, 1, 0], [1, 1, 0], [1, 0, 1]]
25
26   graphs = [gra1, gra2, gra3]
27   print(graph(graphs))
28
```

```
[[0, 7, 6], [7, 0, 5], [6, 5, 0]]
[[0, 9, 7], [9, 0, 6], [7, 6, 0]]
```

1