# BÀI 2: SỰ TƯƠNG ĐỒNG VÀ CÁC KHOẢNG CÁCH

#### I. Mục tiêu:

Sau khi thực hành xong, sinh viên nắm được:

- Khoảng cách giữa các điểm trong tập dữ liệu số sử dụng chuẩn  $L_p$  với  $p=1,2,\infty$ .
- Sự tương đồng của các điểm trong tập dữ liệu phân loại sử dụng: độ đo Overlap và độ đo tần suất xuất hiện ngược.

## II. Tóm tắt lý thuyết:

#### 1. Khoảng cách giữa các điểm trong tập dữ liệu số:

Cho 2 điểm dữ liệu  $\overline{X} = (x_1 \dots x_n)$  và  $\overline{Y} = (y_1 \dots y_n)$ , khoảng cách giữa 2 điểm dữ liệu này dùng chuẩn  $L_p$  được xác định như sau:

$$Dist(\overline{X}, \overline{Y}) = \left(\sum_{i=1}^{n} |x_i - y_i|^p\right)^{1/p}$$

Các trường hợp đặc biệt của chuẩn  ${\cal L}_p$ là

• p = 1 (Manhattan)

$$Dist(\overline{X}, \overline{Y}) = \left(\sum_{i=1}^{n} |x_i - y_i|\right)$$

• p = 2 (Euclidean)

$$Dist(\overline{X}, \overline{Y}) = \left(\sum_{i=1}^{n} |x_i - y_i|^2\right)^{1/2}$$

•  $p = \infty$ 

$$Dist(\overline{X}, \overline{Y}) = \max_{1 \le i \le n} |x_i - y_i|$$

**Ví dụ:** Cho x, y là tọa độ của 4 điểm như sau:

| x | y           |
|---|-------------|
| 3 | 1           |
| 1 | 2           |
| 2 | 0           |
| 2 | 3           |
|   | 3<br>1<br>2 |

Ta có các ma trận khoảng cách như sau:

• Với p=1, ta có

$$Dist(p_1, p_2) = |3 - 1| + |1 - 2| = 2 + 1 = 3$$

$$Dist(p_1, p_3) = |3 - 2| + |1 - 0| = 1 + 1 = 2$$

$$Dist(p_1, p_4) = |3 - 2| + |1 - 3| = 1 + 2 = 3$$

$$Dist(p_2, p_3) = |1 - 2| + |2 - 0| = 1 + 2 = 3$$

$$Dist(p_2, p_4) = |1 - 2| + |2 - 3| = 1 + 1 = 2$$

$$Dist(p_3, p_4) = |2 - 2| + |0 - 3| = 0 + 3 = 3$$

Khi đó, ma trận khoảng cách  $L_1$  giữa 4 điểm này là

|       | $p_1$ | $p_2$ | $p_3$ | $p_4$ |
|-------|-------|-------|-------|-------|
| $p_1$ | 0     | 3     | 2     | 3     |
| $p_2$ | 3     | 0     | 3     | 2     |
| $p_3$ | 2     | 3     | 0     | 3     |
| $p_4$ | 3     | 2     | 3     | 0     |

• Với p=2, ta có

$$Dist(p_1, p_2) = (|3 - 1|^2 + |1 - 2|^2)^{1/2} = \sqrt{5}$$

$$Dist(p_1, p_3) = (|3 - 2|^2 + |1 - 0|^2)^{1/2} = \sqrt{2}$$

$$Dist(p_1, p_4) = (|3 - 2|^2 + |1 - 3|^2)^{1/2} = \sqrt{5}$$

$$Dist(p_2, p_3) = (|1 - 2|^2 + |2 - 0|^2)^{1/2} = \sqrt{5}$$

$$Dist(p_2, p_4) = (|1 - 2|^2 + |2 - 3|^2)^{1/2} = \sqrt{2}$$

$$Dist(p_3, p_4) = (|2 - 2|^2 + |0 - 3|^2)^{1/2} = 3$$

Khi đó, ma trận khoảng cách  $L_2$  (Euclide) giữa 4 điểm này là

|       | $p_1$      | $p_2$      | $p_3$      | $p_4$      |
|-------|------------|------------|------------|------------|
| $p_1$ | 0          | $\sqrt{5}$ | $\sqrt{2}$ | $\sqrt{5}$ |
| $p_2$ | $\sqrt{5}$ | 0          | $\sqrt{5}$ | $\sqrt{2}$ |
| $p_3$ | $\sqrt{2}$ | $\sqrt{5}$ | 0          | 3          |
| $p_4$ | $\sqrt{5}$ | $\sqrt{2}$ | 3          | 0          |

• Với  $p = \infty$ , ta có

$$Dist(p_1, p_2) = \max(|3 - 1|, |1 - 2|) = \max(2, 1) = 2$$
  
 $Dist(p_1, p_3) = \max(|3 - 2|, |1 - 0|) = \max(1, 1) = 1$   
 $Dist(p_1, p_4) = \max(|3 - 2|, |1 - 3|) = \max(1, 2) = 2$   
 $Dist(p_2, p_3) = \max(|1 - 2|, |2 - 0|) = \max(1, 2) = 2$   
 $Dist(p_2, p_4) = \max(|1 - 2|, |2 - 3|) = \max(1, 1) = 1$   
 $Dist(p_3, p_4) = \max(|2 - 2|, |0 - 3|) = \max(0, 3) = 3$ 

Khi đó, ma trận khoảng cách  $L_{\infty}$ giữa 4 điểm này là

|       | $p_1$ | $p_2$ | $p_3$ | $p_4$ |
|-------|-------|-------|-------|-------|
| $p_1$ | 0     | 2     | 1     | 2     |
| $p_2$ | 2     | 0     | 2     | 1     |
| $p_3$ | 1     | 2     | 0     | 3     |
| $p_4$ | 2     | 1     | 3     | 0     |

#### 2. Sự tương đồng giữa các điểm trong tập dữ liệu phân loại:

Cho D là một tập dữ liệu phân loại chứa N đối tượng xác định trên tập d thuộc tính phân loại mà  $A_K$  là thuộc tính thứ i.  $\mathcal{A}_i$  là thuộc tính  $A_i$  lấy  $n_i$  giá trị trong tập D. Xét 2 bản ghi  $\overline{X} = (x_1 \dots x_d)$  và  $\overline{Y} = (y_1 \dots y_d)$ , sự tương đồng đơn giản nhất giữa 2 bản ghi này được xác định như sau

$$Sim(\overline{X}, \overline{Y}) = \sum_{i=1}^{d} w_i S(x_i, y_i)$$

a. Độ đo Overlap:  $S(x_i, y_i)$  là sự tương đồng giữa các giá trị thuộc tính  $x_i, y_i$ 

$$S(x_i, y_i) = \begin{cases} 1 & \text{n\'eu } x_i = y_i \\ 0 & \text{ngược lại} \end{cases}$$

và 
$$w_i = \frac{1}{d}$$
 với  $i = 1 \dots, d$ .

**b.** Độ đo tần suất xuất hiện ngược: Cho  $f_i(x)$  là số lần thuộc tính  $A_i$  lấy giá trị x trong tập dữ liệu D. Nếu  $x \notin A_i$  thì  $f_i(x) = 0$ . Cho  $p_i$  là xác suất của thuộc tính

 $A_i$  lấy giá trị xtrong tập dữ liệu D và được cho bởi công thức

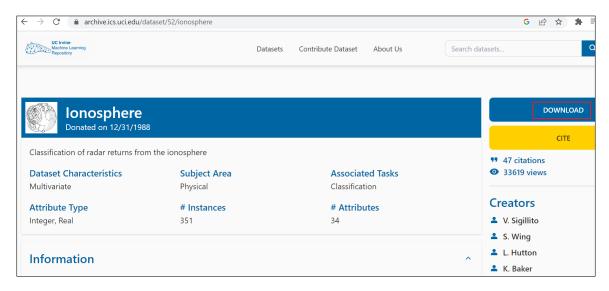
$$p_i(x) = \frac{f_i(x)}{N}$$

$$S(x_i, y_i) = \begin{cases} 1 & \text{n\'eu } x_i = y_i \\ \\ \frac{1}{1 + \log f_i(x_i) \times \log f_i(y_i)} & \text{ngược lại} \end{cases}$$

$$va w_i = \frac{1}{d} v \acute{\alpha} i = 1 \dots, d.$$

### III. Nội dung thực hành:

- 1. Khoảng cách giữa các điểm trong dữ liệu số
  - Download the Ionosphere data set from the UCI Machine Learning Repository (https://archive.ics.uci.edu/dataset/52/ionosphere)



- Đọc dữ liệu từ file "ionosphere.data":

```
import pandas as pd
import numpy as np
df = pd.read_csv('D:\\Huynh\\DataMining_Lab\\data\\tuan2\\ionosphere.data', header=None)
               0.99539 -0.05889
                                  0.85243
                                                0.42267 -0.54487
                                                                   0.18641 -0.45300
               1.00000 -0.18829
                                  0.93035
                                               -0.16626 -0.06288 -0.13738 -0.02447
                                           . . .
               1.00000 -0.03365
                                  1.00000
                                                0.60436 -0.24180
                                                                  0.56045 -0.38238
               1.00000 -0.45161
                                  1.00000
                                                0.25682
                                                        1.00000 -0.32382
                                                                            1.00000
               1.00000 -0.02401
                                  0.94140
                                           . . .
                                               -0.05707 -0.59573 -0.04608
                                                                           -0.65697
               0.83508
                        0.08298
                                                0.86660 -0.10714
                                                                   0.90546 -0.04307
                                           . . .
               0.95113
                                                0.94066 -0.00035
                        0.00419
                                  0.95183
                                                                   0.91483
                                                                           0.04712
  348
            0
               0.94701 -0.00034
                                  0.93207
                                                0.92459 0.00442
                                                                   0.92697 -0.00577
                                           . . .
                                  0.98122
                                                                   0.87403 -0.16243
  349
            0
               0.90608 - 0.01657
                                                0.96022 - 0.03757
  350
              0.84710 0.13533
                                  0.73638
                                                0.75747 -0.06678
                                                                   0.85764 -0.06151
  [351 rows x 35 columns]
```

- Xử lý dữ liệu (bỏ cột cuối):

```
df.pop(df.columns[-1])
print(df)
```

```
... 0.42267 -0.54487 0.18641 -0.45300
             0.99539 -0.05889
                                ... -0.16626 -0.06288 -0.13738 -0.02447
          Λ
             1.00000 -0.18829
                                ... 0.60436 -0.24180 0.56045 -0.38238
... 0.25682 1.00000 -0.32382 1.00000
            1.00000 -0.03365
             1.00000 -0.45161
                                ... -0.05707 -0.59573 -0.04608 -0.65697
            1.00000 -0.02401
                                ...
346
                       0.08298
            0.83508
                                      0.86660 -0.10714
                                                         0.90546 -0.04307
                                . . .
347
             0.95113 0.00419
                                     0.94066 -0.00035
                                                         0.91483 0.04712
          0
             0.94701 -0.00034
                                                         0.92697 -0.00577
348
          0
                                ...
                                      0.92459 0.00442
349
          0
             0.90608 -0.01657
                                      0.96022 -0.03757
                                                         0.87403 -0.16243
             0.84710 0.13533
                                      0.75747 -0.06678
                                                         0.85764 -0.06151
                                 . . .
[351 rows x 34 columns]
```

- Khởi tạo các điểm point<br/>1, point<br/>2, point<br/>3 tương ướng là dòng 0, 1, 2 của array và tính chuẩn  $p=1,2,\infty$ :

```
#array
array = df.values
print(array)
point1 = array[0,:]
point2=array[1,:]
point3 = array[2,:]
dist01_2 = np.linalg.norm(point1 - point2,1)
dist01 3 = np.linalg.norm(point1 - point3,1)
#p=2
dist1 2 = np.linalg.norm(point1 - point2)
dist1_3 = np.linalg.norm(point1 - point3)
#p=inf
dist11 2 = np.linalg.norm(point1 - point2,np.inf)
dist11 3 = np.linalg.norm(point1 - point3,np.inf)
#print results
print(dist1_2)
print(dist1 3)
print(dist0\frac{1}{2})
print(dist01 3)
print(dist11_2)
print(dist11
```

```
0.18641 -0.453
                    0.99539 ... -0.54487
           0.
                            ... -0.06288 -0.13738 -0.02447]
 1.
           0.
                    1.
           0.
[ 1.
                                           0.56045 - 0.38238
                    1.
                             ... -0.2418
[ 1.
           0.
                    0.94701 ... 0.00442
                                           0.92697 -0.00577]
           0.
                    0.90608 ... -0.03757
                                           0.87403 -0.16243]
 1.
                                           0.85764 -0.06151]]
           0.
                    0.8471
                            ... -0.06678
```

```
2.7763589251571923
1.1697276018372824
13.080950000000001
5.35971
1.12221
0.45772
```

- 2. Độ đo tương đồng giữa các điểm trong tập dữ liệu phân loại
  - Download the KDD Cup Network Intrusion Data Set for the UCI Machine Learning Repository

(https://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html)

- Giải nén file "kddcup.data.gz" và đọc file "kddcup.data.conected"

```
<mark>import</mark> pandas <mark>as</mark> pd
import numpy as np
from sklearn.preprocessing import KBinsDiscretizer
from sklearn.preprocessing import StandardScaler
df = pd.read_csv('D:\\Huynh\\DataMining_Lab\\data\\tuan2\\kddcup.data\\kddcup.data.corrected',
                   header = None)
print(df)
                                                                                       40
                                                             0.00
                                                                   0.0
                                                                                0.0
                                                                                      0.0
                     0
                                     SF
                                          215
                                               45076
                                                                         0.00
                         tcp
                              http
                                                                                           normal.
                     0
                         tcp
                              http
                                     SE
                                          162
                                                4528
                                                             0.00
                                                                   0.0
                                                                         0.00
                                                                                0.0
                                                                                      0.0
                                                                                           normal.
                                                       . . .
                         tcp
                     0
                              http
                                     SF
                                          236
                                                1228
                                                             0.00
                                                                   0.0
                                                                         0.00
                                                                                0.0
                                                                                      0.0
                                                                                            normal
                                                       . . .
                     0
                         tcp
                              http
                                     SF
                                          233
                                                2032
                                                             0.00
                                                                   0.0
                                                                         0.00
                                                                                0.0
                                                                                      0.0
                                                                                           normal
                                                       . . .
                              http
                     0
                                     SF
                                          239
                                                 486
                                                             0.00
                                                                   0.0
                                                                         0.00
                                                                                0.0
                                                                                      0.0
                         tcp
                                                       . . .
                                                                                           normal
                                                       . . .
                                          212
          4898426
                              http
                                     SE
                                                2288
                                                             0.05
                                                                   0.0
                                                                         0.01
                                                                                0.0
                                                                                      0.0
                     0
                         tcp
                                                       . . .
                                                                                           normal.
          4898427
                     0
                         tcp
                              http
                                     SF
                                          219
                                                 236
                                                             0.05
                                                                   0.0
                                                                         0.01
                                                                                0.0
                                                                                      0.0
                                                                                           normal
                         tcp
                              http
          4898428
                     0
                                     SF
                                          218
                                                 3610
                                                       . . .
                                                             0.05
                                                                   0.0
                                                                         0.01
                                                                                0.0
                                                                                      0.0
                                                                                           normal.
          4898429
                     0
                         tcp
                              http
                                     SF
                                         219
                                                1234
                                                             0.05
                                                                   0.0
                                                                         0.01
                                                                                0.0
                                                                                      0.0
                                                                                           normal.
                                                       . . .
                                                                   0.0
                                                                                0.0
                                                                                      0.0
          4898430
                                          219
                                                1098
                                                            0.05
                                                                         0.01
                         tcp
                              http
                                                                                           normal.
                                                       . . .
           [4898431 rows x 42 columns]
```

- Chọn các cột có thuộc tính là 'object' (cột 1, 2, 3 và 41)

```
categorical_columns = df.select_dtypes(include=['object']).columns
print("Categorical attributes:", categorical_columns)
categorical_data = df[categorical_columns]
print("Categorical Data:\n", categorical_data)
```

```
Categorical attributes: Int64Index([1, 2, 3, 41], dtype='int64')
Categorical Data:
           1
         tcp
              http
                    SF
                         normal.
         tcp
              http
                    SF
                         normal.
              http
                    SF
                         normal.
         tcp
                    SF
         tcp
              http
                         normal.
                    SF
         tcp
              http
                         normal.
4898426
              http
                     SE
         tcp
                         normal.
4898427
         tcp
              http
                    SF
                         normal.
         tcp
              http
4898428
                     SF
                         normal.
4898429
         tcp
              http
                         normal.
4898430
         tcp
                    SF
                         normal.
              http
[4898431 rows x 4 columns]
```

- Kiểm tra các giá tri giống nhau và loại bỏ những dòng giống nhau này

```
#- Kiểm tra giá trị giống nhau
print("So luong dong giong nhau: ",categorical_data.duplicated().sum())

##loại bỏ các các dòng giống nhau
dfl=categorical_data.drop_duplicates()
print("Du lieu sau khi xoa nhung dong giong nhau: \n",dfl)
```

```
So luong dong giong nhau:
                            4897822
Du lieu sau khi xoa nhung dong giong nhau:
                      2
                                                41
                   http
                                         normal.
1288
         tcp
                   http
                           S2
                                         normal.
1484
         tcp
                   http
                                         normal.
2062
                   smtp
                           SF
                                         normal.
         tcp
              domain_u
2067
         udp
                           SF
                                         normal.
4574130
                   bgp
                                        neptune.
         tcp
                          REJ
4575586
         tcp
                Z39_50
                          REJ
                                        neptune.
4578870
                   smtp
                                        neptune.
         tcp
4579297
         tcp
                 other
                         RSTO
                                        neptune.
4810949
         tcp
                telnet
                         RSTO
                               buffer_overflow.
[609 rows x 4 columns]
```

#### 3. Yêu cầu:

- Viết hàm tính các chuẩn  $p=1,2,\infty$  cho 50 dòng đầu tiên của array trong mục 1.
- Tính các láng giềng gần nhất ở mục 2 sử dụng mục 2 với độ đo Overlap và độ đo tần suất xuất hiện ngược.
- Viết file báo cáo.