# Comp 543 HW6-Outliers

Kai-Po Lin (kl72)

【Code】

```python
import heapq as hq
import numpy as np
import time


def task1(all_data, k, m):
      # The priority queue of outliers
    outliers = list()

    #YOUR CODE HERE!
    outliersDict = dict()
    for idx, i in enumerate(all_data):
        maxPriority = list()
        hq.heapify(maxPriority)
        for jdx, j in enumerate(all_data):
            if idx == jdx:
                continue
            hq.heappush(maxPriority, -np.linalg.norm(j - i))
            if len(maxPriority) > k:
                hq.heappop(maxPriority)
        # Insert idx into with key max(maxPriority)
        outliersDict[-hq.heappop(maxPriority)] = idx

        if len(outliersDict) > m:
            minPriority = list(outliersDict.keys())
            hq.heapify(minPriority)
            outlierKey = hq.heappop(minPriority)
            del outliersDict[outlierKey]

    # Get key by dict() value
    for key, value in outliersDict.items():
        outliers.append((key, value))

    return outliers


def task2(all_data, k, m):
    # Randomly shuffle the data
```

```python
        np.random.shuffle(all_data)

        # The priority queue of outliers
        outliers = list()

        outliersDict = dict()
        minOutlierVal = 0

        for idx, i in enumerate(all_data):
            flag = False
            maxPriority = list()
            hq.heapify(maxPriority)
            for jdx, j in enumerate(all_data):
                if idx == jdx:
                    continue
                hq.heappush(maxPriority, -np.linalg.norm(j - i))
                if len(maxPriority) > k:
                    hq.heappop(maxPriority)
                if len(maxPriority) == k and len(outliersDict) == m and -
hq.nsmallest(1, maxPriority)[0] < minOutlierVal:
                    flag = True
                    break
            if flag:
                continue
            # Insert idx into with key max(maxPriority)
            outliersDict[-hq.heappop(maxPriority)] = idx

            if len(outliersDict) >= m:
                minPriority = list(outliersDict.keys())
                hq.heapify(minPriority)
                minOutlierVal = hq.nsmallest(1, minPriority)[0]

                if len(outliersDict) > m:
                    outlierKey = hq.heappop(minPriority)
                    del outliersDict[outlierKey]

        # Get key by dict() value
        for key, value in outliersDict.items():
            outliers.append((key, value))

        return outliers
```

```python
if __name__ == '__main__':
    # Create the covariance matrix
    covar = np.zeros((100,100))
    np.fill_diagonal(covar, 1)

    # And the mean vector
    mean = np.zeros(100)

    # Create 3000 data points
    all_data = np.random.multivariate_normal(mean, covar, 3000)

    # Now create the 20 outliers
    for i in range(1, 20):
        mean.fill(i)
        outlier_data = np.random.multivariate_normal(mean, covar, i)
        all_data = np.concatenate((all_data, outlier_data))

    # k for kNN detection
    k = 10

    # The number of outliers to return
    m = 5

    # YOUR CODE HERE!

    # Task 1
    # Start the timer
    start_time = time.time()

    outliers = task1(all_data, k, m)

    print("Task1:")
    print("--- %s seconds ---" % (time.time() - start_time))

    # Print the outliers...
    for outlier in outliers:
        print(outlier)


    # Task 2
    # Start the timer
    start_time = time.time()
```

```
outliers = task2(all_data, k, m)


print("\nTask2:")
print("--- %s seconds ---" % (time.time() - start_time))


# Print the outliers...
for outlier in outliers:
    print(outlier)
```

【Result】=> Task 1: About 89.08 seconds, Task 2: About 4.13 seconds

```
PS C:\Users\KB\Desktop\Rice\Courses\Tools & Models for DS\HW\HW6-Outliers> python .\outlier.py
Task1:
--- 89.079591274261147 seconds ---
(21.99866787819994, 3002)
(24.237764820644813, 3003)
(24.21031384711911, 3004)
(23.133842954427642, 3005)
(20.280561699431107, 3007)

Task2:
--- 4.1308159828186035 seconds ---
(24.237764820644813, 563)
(24.21031384711911, 1231)
(21.99866787819994, 2319)
(23.133842954427642, 2440)
(20.280561699431107, 2580)
```