

# 情報工学科 レポート

実験演習記録			判定・指示		
	年月日時	共同作業者			
1	2019/06/17				
2					
3					
4					
レポート提出記録					
	提出年月日	期限年月日			
初	2019/07/01	2019/07/01			
再					
科目名		テーマ担当教員	学年	学期	単位
情報工学実験2		中川先生	3	前期	2
テーマ番号	テーマ名	学籍番号 氏 名			
	パターン認識	16268062 SALIC ERTUGRUL			

## 1. はじめに

本レポートでは、情報工学実験 2 の第 8 回目のレポートとして、パターン認識についての課題に取り組み、作成するプログラムの設計や実装について述べる。

## 2. 目的

本実験では、パターン認識でよく用いられる「クラスタリング」を、簡単なプログラムの作成を通して理解することを目的とする。

## 3. 実験方法

本実験では、全ての課題を EDEN において、「Python 言語」で行う。

この実験では、k-means 法によってクラスタリングを行う。そのため、まず、いくつかのクラスターを持つ 2 次元特徴のパターンを考える。そのパターンは <https://archive.ics.uci.edu/ml/datasets/iris> (閲覧日: 2019 年 7 月 1 日) からダウンロードしたデータを 2 次元に変えたものである。(レポート提出時に添付する) データは txt ファイル (iris4d.txt) に保存されているため、プログラムでは最初にファイルからデータを読み込む処理を行う。そのコードを次に示す。

```
1. import numpy as np
2. X = np.loadtxt("./iris4d.txt")
3. (d, n) = X.shape # 行数と列数を取得
```

クラスタリングを行う際は mahalanobis 距離を用いる。mahalanobis 距離は式 (1) で定義できる。

$$d(x, y) = (x - y)^T \Sigma (x - y) \quad (1)$$

式 (1) を用いて、mahalanobis 距離を求める関数を定義する。共分散行列はデータ全体の共分散行列だと仮定する。そのコードを次に示す。

```
1. def dis(x, y, cov): # Mahalanobis 距離の計算
2.     dist = np.dot(np.dot((x-y).T, cov), (x-y))
3.     return dist
```

k-means 法を以下の STEP で行う。

- STEP-1: 有限パターン集合  $W$  をいくつか分割するかを決める. これを  $N$  とし, 指定できるようにする.
- STEP-2:  $N$  分割する時, それぞれ代表パターンと思われる  $c_1, c_2, c_3, \dots, c_N$  を適当に選ぶ.
- STEP-3: 全てのパターン  $x$  について, もっとも近い代表パターン  $c_i$  を求め,  $x$  が属する  $S_i$  を決める.
- STEP-4: 各  $S_i$  について, セントロイドを求め, 代表パターンをそれに更新する.
- STEP-5: もし全ての代表パターンに変化がなければ終了.
- STEP-6: STEP-3 に戻る.

STEP-1 ではユーザに  $N$  を指定してもらう. そのコードは次のように書ける.

```
1. # STEP-1
2. print("N:")
3. N = int(input())
```

また, STEP-2 ではセントロイドの初期値を適当に定義する. 今回はデータの最初の  $N$  個のパターンを代表パターンとして決めておく. そのコードを次に示す.

```
1. # STEP-2
2. c = np.zeros((d, N))
3. for i in range(N):
4.     c[:, i] = X[:, i]
5. print("セントロイドの初期値:")
6. print(c)
```

STEP-3 では, 全てのパターン  $x$  について, もっとも近い代表パターン  $c_i$  を求め,  $x$  が属する  $S_i$  を決める. そこでマハラノビス距離を用いて,  $x$  が一番近い代表パターン  $c_i$  を求める. そのコードを次に示す.

```
1. S = []
2. for i in range(N):
3.     S.append([]) # N個の行列を要素とするリスト
4. d_mhlnbs = [0]*N
5. for i in range(n):
6.     for j in range(N):
7.         d_mhlnbs[j] = dis(X[:, i], c[:, j], cov) #マハラノビス距離
8.     min_d = d_mhlnbs.index(min(d_mhlnbs)) # 一番近い代表のクラスを取得
9.     S[min_d].append(X[:, i]) # 一番近い代表のクラスに入れる
```

STEP-4 では, 各  $S_i$  について, セントロイドを求め, 代表パターンをそれに更新する. セントロイドを計算する際は,  $S$  に所属している  $x$  から  $S$  の各要素への距離の合計を求め, その合計値が一番小さくなるような  $x$  を計算する. そのコードを次に示す.

```

1. # セントロイドを更新
2. for i in range(N):
3.     for j in range((len(S[i]))):
4.         summ = 0
5.         for k in range((len(S[i]))):
6.             summ = summ + dis(S[i][j], S[i][k], cov)
7.         if j == 0:
8.             min_value = sum
9.             min_index = j
10.        elif summ < min_value:
11.            min_value = sum
12.            min_index = j
13.        c[:, i] = S[i][min_index]
14.    print("新しいセントロイド:")
15.    print(c)

```

STEP-5 では、全ての代表パターンに変化がなければプログラムを終了する。そこで前のセントロイドを保存しておく **control** 変数を用いる。もし更新したセントロイドが **control** と同じであればプログラムを終了する。でなければまた処理を繰り返す。そのコードを次に示す。

```

1. # 前のセントロイドと更新したものが同じであれば終了
2. for i in range(d):
3.     for j in range(N):
4.         if c[i, j] != control[i, j]:
5.             flag = 1

```

もし **control**（前のセントロイド）と更新したセントロイドが同じであれば **flag** が 0 のままになり、プログラムが終了する。一方、同じでなければ **flag** が経ち、処理が繰り返される。つまり STEP-3 に戻る。

また、クラスタリングしたデータをグラフで表示したい。そのために次にコードを用いる。

```

1. # グラフの用意
    fig = plt.figure()
2. ax = fig.add_subplot(1, 1, 1)
3. color_list = ["r", "g", "b", "c", "m", "y", "k", "w"]
4. for i in range(N):
5.     for j in range(len(S[i])):
6.         ax.scatter(S[i][j][0], S[i][j][1], c=color_list[i %
            len(color_list)], marker='.', label='group1')
7. fig.show()

```

## 4. 実装結果

本節では、節3で説明したプログラムの実装結果について述べる。

プログラムを実装し、 $N=3$ にした結果を図1に示す。また、その時出力されたグラフを図2に示す。

```
Run: kadai1 x
"C:\Users\Katre Jp\AppData\Local\Prog
N:
3
セントロイドの初期値 :
[[5.1 4.9 4.7]
 [3.5 3. 3.2]]
新しいセントロイド :
[[6.1 5.1 4.6]
 [3. 2.5 3.1]]
新しいセントロイド :
[[6.4 5.4 4.8]
 [2.9 3. 3.4]]
新しいセントロイド :
[[6.6 5.6 4.8]
 [3. 3. 3.4]]
新しいセントロイド :
[[6.7 5.6 4.9]
 [3. 3. 3.1]]
新しいセントロイド :
[[6.7 5.7 4.9]
 [3. 3. 3.1]]
新しいセントロイド :
[[6.7 5.7 4.9]
 [3. 3. 3.1]]
Process finished with exit code 0
```

図1 クラスタリング実装結果

図1ではセントロイドの状況や変化が見られる。図2ではクラスタリングの結果が確認できる。

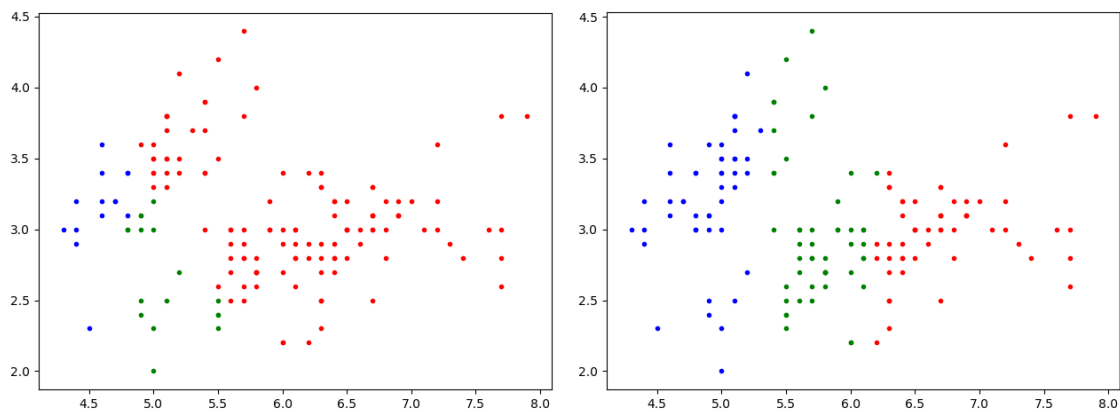


図2  $N=3$ の時のクラスタリング結果

また、N=5 の時のクラスタリングを図 3 に示す.

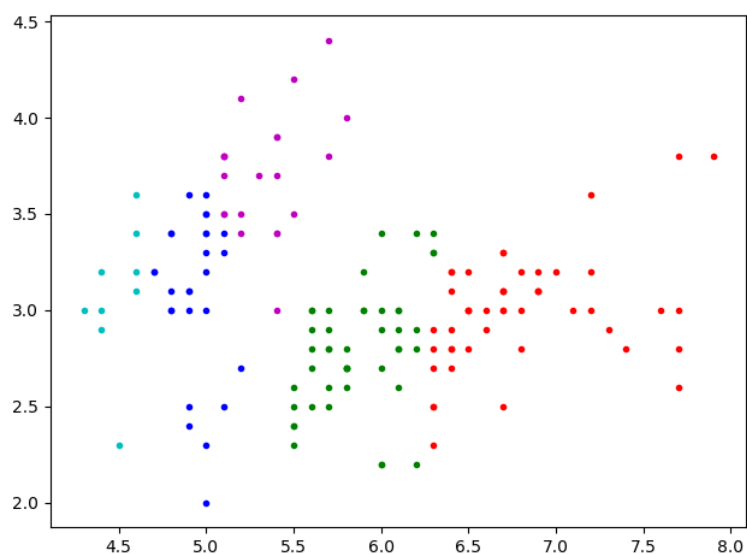


図 3 N=5 の時のクラスタリング結果

最後に、N=3 とし、セントロイドの初期値を変えてみた時の結果を出してみる. 今まででは、セントロイドの初期値をデータの最初の N 個の要素としていたが、それを最後の N 個の要素としてみる. その時の結果を図 4 と図 5 に示す.

```
Run: kadai1 x
N:
3
セントロイドの初期値 :
[[4.7 4.9 5.1]
 [3.2 3. 3.5]]
新しいセントロイド :
[[4.6 5.1 6.1]
 [3.1 2.5 3. ]]
新しいセントロイド :
[[4.8 5.4 6.4]
 [3.4 3. 2.9]]
新しいセントロイド :
[[4.8 5.6 6.6]
 [3.4 3. 3. ]]
新しいセントロイド :
[[4.9 5.7 6.7]
 [3.1 3. 3. ]]
新しいセントロイド :
[[4.9 5.7 6.7]
 [3.1 3. 3. ]]
Process finished with exit code 0
```

図 4 セントロイドの初期値を変えた時の結果 (N=3)

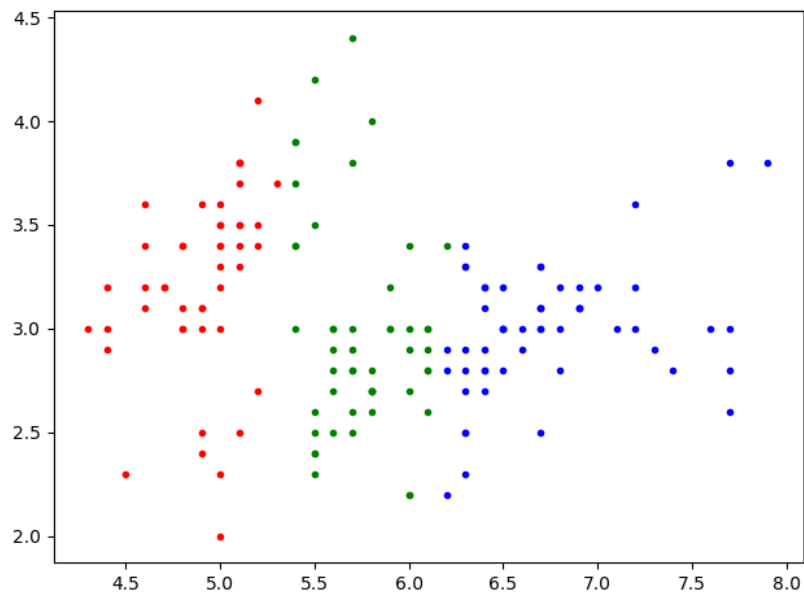


図5 セントロイドの初期値を変えた時の結果 (N=3)

## 5. 考察

本節では、節4の結果を基に、実験に関して考察を行うことを目的とする。

本実験では、k-means法を用いてクラスタリングを行った。図2(右)より、クラスタリングのプログラムが正常に動いていると考えられる。そこでセントロイドの初期値をパターン集合から適当に選んだので、図2(左)が、まだ完璧な結果になっていないことも分かる。図1ではセントロイドの変化が見られる。セントロイドとしてプログラムが適切としているパターンの行列は

$$\begin{pmatrix} 6.7 & 5.7 & 4.9 \\ 3.0 & 3.0 & 3.1 \end{pmatrix}$$

だと分かる。その行列に格納されているセントロイドの値が、セントロイドの初期値を変えてプログラムを実行した時の結果と、つまり最終セントロイドの値と同じであることは、注意すべきだと考えられる。したがって、プログラムは正常に動いていることがこの点でも確認できている。またそれは、図2(右)と図5のグラフが同じであることから確認できる。

本実験では、またN=5とし、クラスタリングを行ってみた。その結果が図3に示されている。図3より、クラス数が変わってもクラスタリングができることが確認できる。

本実験で使用したデータは3種類の花のデータである。これから、そのデータの本当のクラスについて紹介する。データの本当のクラスは図6に示す。

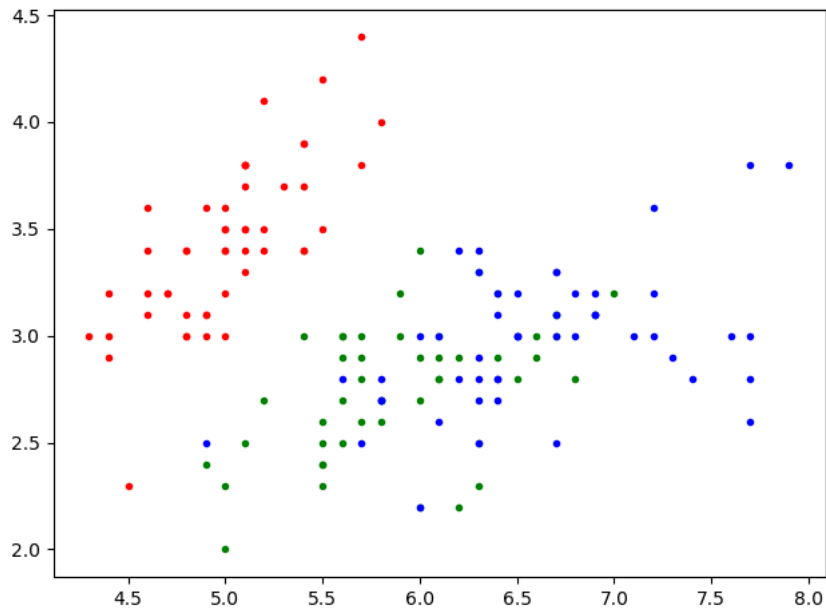


図6 本実験で用いたデータの本当のクラスのグラフ

図5のN=3の時のプログラム実装結果と図6の本当のクラスを比較してみたところ、プログラムが行ったクラスタリングははっきりとクラスを分けることができていることに気づいた。データの本当のクラスでは(図6)緑色のグループと青色のグループがはっきりと分かれていないことがわかるが、図5でははっきりと分けられている。よって、プログラムの分類精度が高いと思われる、しかしそれにしても、実際と誤差が生じていることがわかる。

## 6. おわりに

本実験ではk-means法を用いてクラスタリングを行った。それは情報工学科の学生である私にとって、とても貴重な経験・知識であった。なぜそうであるかを、このレポートを書く時に起こった出来事で説明させて頂きたい。

本レポートを書いているところに、東京の他の大学で国際関係について学んでいる友達が来た。私がk-means法についてプログラミングをしているのを見て、「あっ！、これ最近使った方法。情報工学科の人は偉そうにIT系の人だけやっているように思ってるかもしれないけど、僕たちもこういうのたまにやってるよ」と冗談付きで言った。私は同じ学年の文系の友達が、k-means法を使って分類等を行ったことにびっくりしつつ、理系しかも情報系の勉強をしている私が今までk-means法の使い方を良く知らなかったことを恥ずかしく思った。本実験で非常に重要な知識を身についたと思う。(ありがとうございます..)

## 7. 参考文献

- [1] 中川 正樹, 東京農工大学, 情報工学科, 情報工学実験2, 2019年度第8回講義資料
- [2] <<https://archive.ics.uci.edu/ml/datasets/iris>> (閲覧日: 2019年7月1日)
- [3] <[https://octave.org/doc/v4.2.1/Two\\_002dDimensional-Plots.html](https://octave.org/doc/v4.2.1/Two_002dDimensional-Plots.html)>, 閲覧日: 2019年6月30日
- [4] <<https://pythondatascience.plavox.info/matplotlib/pycharm>>, 閲覧日: 2019年6月30日