

```
+ コード + テキスト
RAM
ディスク
k近傍法
```

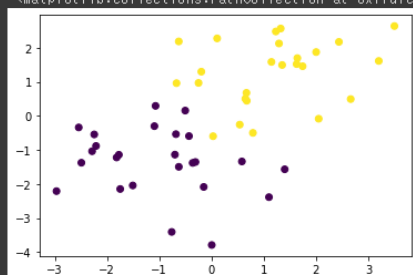
```
1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from scipy import stats
```

訓練データ生成

```
[16] 1 def gen_data():
2     #50の要素を持つベクトルを2列に変換して各要素から-1
3     x0 = np.random.normal(size=50).reshape(-1,2) - 1
4     #50の要素を持つベクトルを2列に変換して各要素から+1
5     x1 = np.random.normal(size=50).reshape(-1,2) + 1.
6
7     x_train = np.concatenate([x0,x1])
8     y_train = np.concatenate([np.zeros(25),np.ones(25)]).astype(np.int)
9     return x_train,y_train
10
11 # r1,r2 = gen_data()
12 # print(np.mean(r1[0]))
13 # print(np.mean(r1[1]))
14
```

```
[17] 1 X_train, ys_train = gen_data()
2 print(np.mean(X_train[:,0]))
3 print(np.mean(X_train[:,1]))
4 # print(X_train)
5 plt.scatter(X_train[:,0], X_train[:,1], c=ys_train)
```

```
0.0751073106537703
-0.04581431581136874
<matplotlib.collections.PathCollection at 0x7fd1cfff36d0>
```



学習 : k近傍法において、「学習」相当の処理はなし

予測するデータ点との、距離が最も近いk個の、訓練データのラベルの最頻値を割り当てる。

```
[18] 1 #引数1,2の「差の二乗」の合計を返却
2 def distance(x1, x2):
3     return np.sum((x1 - x2)**2, axis=1)
4
5 #k近傍法予測
6 def knn_predict(n_neighbors, x_train, y_train, X_test):
```

0秒 完了時間: 21:50

```
3     return np.sum((x1 - x2)**2, axis=1)
4
5 #k近傍法予測
6 def knn_predict(n_neighbors, x_train, y_train, X_test):
7     ##### np.empty() #####
8     #用途: 未初期化の配列を生成する。
9     #第1引数: shape
10    #第2引数: 要素のタイプ(初期値float)
11    y_pred = np.empty(len(X_test), dtype=y_train.dtype) #X_testと同じサイズの行列を初期化生成
12    j = 0
13    for i, x in enumerate(X_test):
14        distances = distance(x, X_train)
15        nearest_index = np.argsort(distances)[:n_neighbors]
16        #scipy.stats.mode: 最頻値 (モード)
17        mode, uu = stats.mode(y_train[nearest_index])
18        y_pred[i] = mode
19        # if i>5000 and i<6000:
20        #     print(y_pred[i])
21    return y_pred
22
```

```

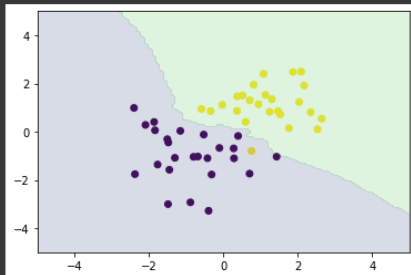
23
24 def plt_result(x_train, y_train, y_pred):
25     xx0, xx1 = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
26     xx = np.array([xx0, xx1]).reshape(2, -1).T
27
28     plt.scatter(x_train[:,0], x_train[:,1], c=y_train)
29     plt.contourf(xx0, xx1, y_pred.reshape(100,100).astype(dtype=np.float), alpha=0.2, levels=np.linspace(0, 1, 3))
30
31

```

```

[5] 1 n_neighbors = 3
2
3 xx0, xx1 = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
4 X_test = np.array([xx0, xx1]).reshape(2, -1).T
5
6 # print(np.array([xx0, xx1]))
7 # print(np.array([xx0, xx1]).flatten())
8 # print(np.array([xx0, xx1]).reshape(2,-1))
9 y_pred = knn_predict(n_neighbors, X_train, ys_train, X_test)
10 plt_result(X_train, ys_train, y_pred)
11

```



```

[6] 1 #skilにて実装

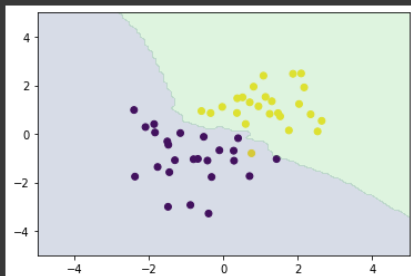
```

+ コード + テキスト

```

2 from sklearn.neighbors import KNeighborsClassifier
3 knn = KNeighborsClassifier(n_neighbors).fit(X_train, ys_train)
4 xx = np.array([xx0, xx1]).reshape(2, -1).T
5 # print(xx)
6 # xx = np.array([xx0, xx1]).reshape(-1, 2)
7 # print(xx)
8 plt_result(X_train, ys_train, knn.predict(xx))
9

```



▼ k平均クラスタリング(k-means)

```

[7] 1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4

```

データ作成

```

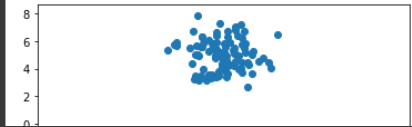
[8] 1 def gen_data():
2     x1 = np.random.normal(size=(100, 2)) + np.array([-5, -5])
3     x2 = np.random.normal(size=(100, 2)) + np.array([5, -5])
4     x3 = np.random.normal(size=(100, 2)) + np.array([0, 5])
5     #np.vstack:列方向に行列を結合 (列数がある必要ある。)
6     #np.hstack:行方向に行列を結合 (行数がある必要ある。)
7     return np.vstack((x1,x2,x3))
8
9 # a = np.array([5, -5])
10 # b = np.random.normal(size=(100, 2))
11 # c = a + b
12
13 # print(a)
14 # print(b[:5])

```

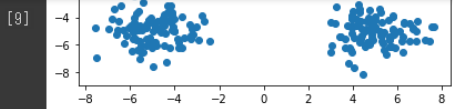
```
15 # print(c[:5])
16 # print(gen_data().shape)
17
```

```
[9] 1 #データ作成
    2 X_train = gen_data()
    3 #データ描画
    4 plt.scatter(X_train[:, 0], X_train[:, 1])
    5
```

<matplotlib.collections.PathCollection at 0x7fd1d2159090>



+ コード + テキスト



✓ RAM
ディスク

▼ 学習(クラスタ中心の計算)

k-meansアルゴリズムは以下の通り。

- 1) 各クラスタ中心の初期値を設定する。
- 2) 各データ点に対して、各クラスタ中心との距離を計算し、最も距離が近いクラスタを割り当てる。
- 3) 各クラスタの平均ベクトル (中心) を計算する。
- 4) 収束するまで 2, 3 の処理を繰り返す

```
[10] 1 def distance(x1, x2):
    2     #axis=1は集計する方向を示す。
    3     return np.sum((x1-x2)**2, axis=1)
    4
    5 # a1 = np.array([[1,2,3],[4,5,6],[1,2,3],[4,5,6]])
    6 # a2 = np.array([[4,5,6],[1,2,3],[1,2,3],[4,5,6]])
    7 # np.sum((a1 - a2)**2, axis=0)
    8
    9 n_clusters = 3
   10 iter_max = 100
   11
   12 #1) 各クラスタ中心の初期値を設定する。
   13 centers = X_train[np.random.choice(len(X_train), n_clusters, replace = False)]
   14 print("クラスタの初回中心(ランダム):")
   15 print(centers)
   16 print()
   17 for cnt in range(iter_max):
   18     #前ループの各クラスタ中心を保存したうえで、centerの計算を開始
   19     prev_centers = np.copy(centers)
   20     D = np.zeros((len(X_train), n_clusters))
   21
   22     #2-1) 【中心再計算】各データ点に対して、各クラスタ中心との距離を再計算。
   23     #Dは各列に「3中心からの距離」を保持する。
   24     for i, x in enumerate(X_train):
   25         D[i] = distance(x, centers)
   26
   27     #各データ点に、最も距離が近いクラスタを割り当て
   28     #argmin()は配列の最小値のインデックスを取得するもの
   29     #「D(距離配列)」から中心から最小距離を持つインデックスを選択し、リストとして取り出す。(0 or 1 or 2)
   30     #2-2) 【再分類】最も距離が近いクラスタを割り当てる。
   31     cluster_index = np.argmin(D, axis=1)
   32
   33     #各クラスタの中心を計算
   34     #3) 各クラスタの平均ベクトル (中心) を再計算する。
   35     for k in range(n_clusters):
   36         index_k = cluster_index == k
   37         centers[k] = np.mean(X_train[index_k], axis=0)
   38
   39     #ループ終了判定
   40     #前ループのクラスタ中心と今回ループのクラスタ中心が同じであれば学習終了
   41     if np.allclose(prev_centers, centers):
   42         # print(cluster_index)
   43         # print(D)
   44         print(index_k)
   45
   46     # print("最終ループ時の各点距離D:")
   47     # print(D)
```

✓ RAM
ディスク

+ コード + テキスト

```
[10] 46 # print("最終ループ時の各点距離D:")
    47 # print(D)
```