

# SML HW6

29-176004 奥村 恭平<sup>\*</sup>

June 9, 2018

## 宿題 1: カーネル密度推定法の実装

ガウスカーネルに対するカーネル密度推定法を実装した。(言語は python.) バンド幅は尤度交差確認により決定 ( $b^* := 0.1$ ) した. 以下はシミュレーション結果の図. 上がオリジナルのサンプルのヒストグラム, 下が推定された密度関数. うまく推定できていることがわかる.

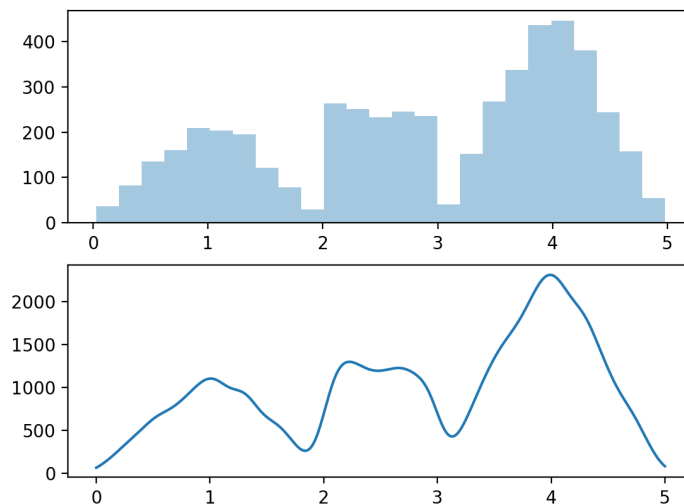


Figure 1: シミュレーション実行結果

```
1 import numpy as np
2 import scipy as sp
3 from numpy.random import randn, rand
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from scipy.stats import norm
7
8
9 def myrand(n=5000):
10     x = np.zeros(n)
11     u = rand(n)
12     flag = (0 <= u) * (u < 1/8)
13     x[flag] = np.sqrt(8*u[flag])
14     flag = (1/8 <= u) * (u < 1/4)
15     x[flag] = 2 - np.sqrt(2 - 8*u[flag])
16     flag = (1/4 <= u) * (u < 1/2)
17     x[flag] = 1 + 4*u[flag]
```

<sup>\*</sup>E-mail: kyohei.okumura@gmail.com

<sup>†</sup>東京大学大学院 経済学研究科 M2

```

18     flag = (1/2 <= u) * (u < 3/4)
19     x[flag] = 3 + np.sqrt(4*u[flag] - 2)
20     flag = (3/4 <= u) * (u <= 1)
21     x[flag] = 5 - np.sqrt(4 - 4*u[flag])
22     return x
23
24
25 def gkernel_est(sample, x=np.linspace(0, 5, 501), bandwidth=0.1):
26     pxh = np.zeros_like(x)
27     n = sample.shape[0]
28     for i in range(n):
29         pxh = pxh + norm.pdf(x, loc=sample[i], scale=bandwidth)
30     return x, pxh
31
32
33 def cross_validation(sample, n_split=5, params=[0.01, 0.1, 0.5]):
34     n_params = len(params)
35     likelihoods = np.zeros(n_params)
36     group = np.split(sample, n_split)
37     for j in range(n_params):
38         for i in range(n_split):
39             if i==0:
40                 sample_temp = np.hstack(group[i+1:][0])
41             elif i==n_split-1:
42                 sample_temp = np.hstack(group[0:i])
43             else:
44                 sample_temp = np.hstack([np.hstack(group[0:i]), group[i
45                                         +1:][0]])
46             _, pxh = gkernel_est(sample_temp, group[i], bandwidth=params[j])
47             likelihoods[j] += np.sum(np.log(pxh))
48         opt_param = params[np.argmax(likelihoods)]
49         #print(likelihoods)
50     return opt_param
51
52 if __name__ == '__main__':
53     np.random.seed(1)
54     sample = myrand()
55     opt_b = cross_validation(sample=sample)
56     x, pxh = gkernel_est(sample=sample, bandwidth=opt_b)
57
58     # plot original samples
59     fig = plt.figure()
60     ax1 = fig.add_subplot(2,1,1)
61     sns.distplot(sample, kde=False, rug=False, bins=25)
62     None
63
64     # plot estimated distribution
65     ax2 = fig.add_subplot(2,1,2)
66     ax2.plot(x, pxh)
67     plt.show()

```

## 宿題 2: kNN 法による文字識別

kNN 法による手書き文字認識を行った。言語は python。最近傍数  $k$  は交差確認により決定した。(正答率の計算部分のみ、scikit-learn のモジュールを用いた。) 訓練データを用いて  $k \in \{1, 2, 3, 4, 5, 10, 20\}$  について交差確認を行ったところ、 $k^* = 1$  が最適になった。(交差確認におけるそれぞれの平均正答率は、 $[0.9676, 0.9676, 0.9674, 0.9688, 0.965, 0.959, 0.9454]$  となった。)  $k := 1$  としてテストデータを識別したところ、正答率は、0.965 となった。

```
1 import numpy as np
2 from collections import Counter
3 from sklearn.metrics import accuracy_score
4
5
6 class kNN(object):
7     def __init__(self, k=1):
8         self._train_data = None
9         self._target_data = None
10        self._k = k
11
12
13    def fit(self, train_data, target_data):
14        self._train_data = train_data
15        self._target_data = target_data
16
17
18    def predict(self, x):
19        distances = np.array([np.linalg.norm(p - x) for p in self._train_data])
20        nearest_indices = distances.argsort()[:self._k]
21        nearest_labels = self._target_data[nearest_indices]
22        c = Counter(nearest_labels)
23        return c.most_common(1)[0][0]
24
25
26 def load_train_data():
27     for i in range(10):
28         if i==0:
29             train_feature = np.loadtxt('data/digit_train{}.csv'.format(i),
30                                         delimiter=',')
31             train_label = np.array([i]*train_feature.shape[0])
32         else:
33             temp_feature = np.loadtxt('data/digit_train{}.csv'.format(i),
34                                       delimiter=',')
35             train_feature = np.vstack([train_feature, temp_feature])
36             temp_label = np.array([i]*temp_feature.shape[0])
37             train_label = np.hstack([train_label, temp_label])
38
39     return train_feature, train_label
40
41 def load_test_data():
42     for i in range(10):
43         if i==0:
44             test_feature = np.loadtxt('data/digit_test{}.csv'.format(i),
45                                       delimiter=',')
46             test_label = np.array([i]*test_feature.shape[0])
47         else:
48             temp_feature = np.loadtxt('data/digit_test{}.csv'.format(i),
49                                       delimiter=',')
50             test_feature = np.vstack([test_feature, temp_feature])
51             temp_label = np.array([i]*temp_feature.shape[0])
52             test_label = np.hstack([test_label, temp_label])
53
54     return test_feature, test_label
```

```

53
54 def calc_accuracy(train_feature, train_label, test_feature, test_label, k=1)
    :
55     model = kNN(k)
56     model.fit(train_feature, train_label)
57     predicted_labels = []
58     for feature in test_feature:
59         predicted_label = model.predict(feature)
60         predicted_labels.append(predicted_label)
61     return accuracy_score(test_label, predicted_labels)
62
63
64 def load_train_data_cv(n_split=5):
65     for i in range(10):
66         if i==0:
67             train_feature = np.loadtxt('data/digit_train{}.csv'.format(i),
68                                     delimiter=',')
69             train_label = np.array([i]*train_feature.shape[0])
70             group_feature = np.split(train_feature, n_split)
71             group_label = np.split(train_label, n_split)
72         else:
73             temp_feature = np.loadtxt('data/digit_train{}.csv'.format(i),
74                                     delimiter=',')
75             temp_group_feature = np.split(temp_feature, n_split)
76             temp_label = np.array([i]*temp_feature.shape[0])
77             temp_group_label = np.split(temp_label, n_split)
78
79             for m in range(n_split):
80                 group_feature[m] = np.vstack([group_feature[m],
81                                             temp_group_feature[m]])
82                 group_label[m] = np.hstack([group_label[m], temp_group_label
83                                             [m]])
84
85     return group_feature, group_label
86
87
88 def cross_validation(n_split=5, params=[1,2,3,4,5,10,20]):
89     n_params = len(params)
90     score_list = np.zeros(n_params)
91     group_feature, group_label = load_train_data_cv(n_split)
92
93     for j in range(n_params):
94         for i in range(n_split):
95             temp_group_feature = group_feature.copy()
96             temp_test_feature = temp_group_feature.pop(i)
97             temp_train_feature = np.vstack(temp_group_feature)
98
99             temp_group_label = group_label.copy()
100            temp_test_label = temp_group_label.pop(i)
101            temp_train_label = np.hstack(temp_group_label)
102
103            score_list[j] += calc_accuracy(temp_train_feature,
104                                        temp_train_label, temp_test_feature, temp_test_label, k=
105                                        params[j])/n_split
106
107     opt_param = params[np.argmax(score_list)]
108     print(score_list)
109     return opt_param
110
111
112 def main():
113     k_opt = cross_validation(n_split=5, params=[1,2,3,4,5,10,20])
114     train_feature, train_label = load_train_data()
115     test_feature, test_label = load_test_data()
116     score = calc_accuracy(train_feature, train_label, test_feature,
117                           test_label, k=k_opt)

```

```
111     print(score)
112
113
114 if __name__ == '__main__':
115     main()
```