# 統計的機械学習
# 第六回　レポート

37-196360　森田涼介

2019 年 5 月 27 日

# 宿題 1

　ガウスカーネルに対するカーネル密度推定法を行う。バンド幅は尤度交差確認法によって決定する。

　結果を以下の表 4 と図 1 に示した。表 4 より，最も尤度の平均の大きいバンド幅は 0.1 となった。なお，プログラムは 4 ページの Listing 1 に示した。

表 1: バンド幅とそれに対応する LCV の値

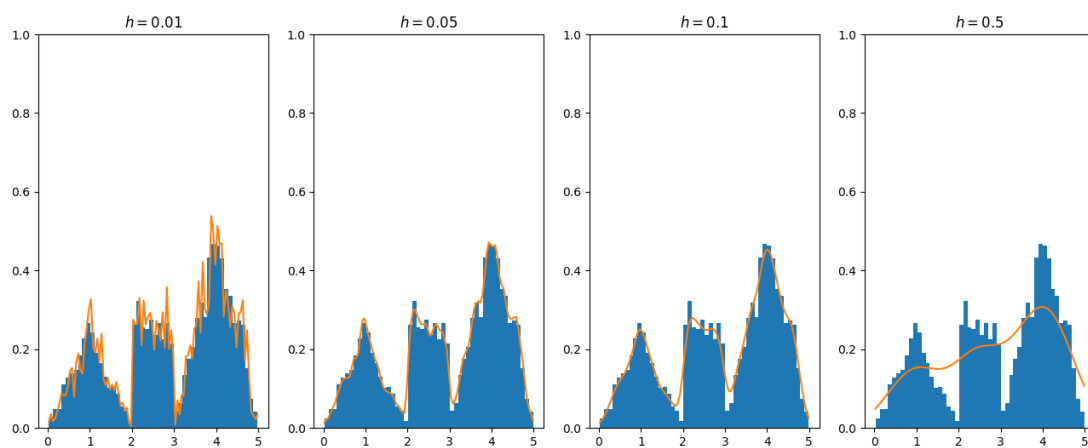| $h$ | 0.01 | 0.05 | 0.10 | 0.50 |
|---|---|---|---|---|
| LCV | $-6344$ | $-4009$ | $-3938$ | $-4219$ |



図 1: 各バンド幅に対するヒストグラムと確率密度

# 宿題 2

　最近傍識別器によって 0–9 までの 10 クラスの手書き文字認識を行う。近傍数 $k$ は識別誤差に関する交差確認により決定する。

　以下に結果を示す。近傍数 $k$ の候補として，1–10 を考えた。各 $k$ についての，交差確認法で求めた正解数の平均値を表 2 に示す。これより，最も高い正解数を与える $k$ を選ぶと，$k = 3$ となる。$k = 3$ に対する混同行列は表 3 のようになり，また，各カテゴリごとの正解率等は表 4 のようになった。

　プログラムは 8 ページの Listing 2 に示した。

表 2: 各 $k$ についての，交差確認法で求めた正解数の平均値

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 484.7 | 483.2 | 485.1 | 484.8 | 483.4 | 482.7 | 481.6 | 481.2 | 480.0 | 479.7 |

表 3: 混同行列

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 198 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 199 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 195 | 0 | 0 | 0 | 0 | 2 | 3 | 0 |
| 3 | 0 | 0 | 0 | 190 | 0 | 4 | 0 | 1 | 4 | 1 |
| 4 | 0 | 1 | 0 | 0 | 189 | 0 | 3 | 0 | 0 | 7 |
| 5 | 2 | 0 | 3 | 3 | 1 | 186 | 0 | 0 | 1 | 4 |
| 6 | 1 | 0 | 2 | 0 | 0 | 0 | 197 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 0 | 4 | 0 | 0 | 191 | 0 | 4 |
| 8 | 2 | 0 | 2 | 3 | 0 | 3 | 0 | 0 | 187 | 3 |
| 9 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 198 |

表 4: 各カテゴリごとの結果

| Category | #Data | #Correct | Accuracy |
|---|---|---|---|
| 0 | 200 | 198 | 0.990 |
| 1 | 200 | 199 | 0.995 |
| 2 | 200 | 195 | 0.975 |
| 3 | 200 | 190 | 0.950 |
| 4 | 200 | 189 | 0.945 |
| 5 | 200 | 186 | 0.930 |
| 6 | 200 | 197 | 0.985 |
| 7 | 200 | 191 | 0.955 |
| 8 | 200 | 187 | 0.935 |
| 9 | 200 | 198 | 0.990 |
| All | 2,000 | 1,930 | 0.965 |

# プログラム

実行環境と用いた言語・ライブラリを以下の表 5 に示す。

<div align="center">

表 5: プログラムの実行環境

| | | |
|---|---|---|
| OS | : | Microsoft Windows 10 Pro (64bit) |
| CPU | : | Intel(R) Core(TM) i5-4300U |
| RAM | : | 4.00 GB |
| 使用言語 | : | Python3.6 |
| 可視化 | : | matplotlib ライブラリ |

</div>

Listings 1: `assignment1.py`

```python
# -*- coding: utf-8 -*-


import numpy as np
import matplotlib.pyplot as plt


def generate_data(n=3000):
    x = np.zeros(n)
    u = np.random.rand(n)
    index1 = np.where((0 <= u) & (u < 1 / 8))
    x[index1] = np.sqrt(8 * u[index1])
    index2 = np.where((1 / 8 <= u) & (u < 1 / 4))
    x[index2] = 2 - np.sqrt(2 - 8 * u[index2])
    index3 = np.where((1 / 4 <= u) & (u < 1 / 2))
    x[index3] = 1 + 4 * u[index3]
    index4 = np.where((1 / 2 <= u) & (u < 3 / 4))
    x[index4] = 3 + np.sqrt(4 * u[index4] - 2)
    index5 = np.where((3 / 4 <= u) & (u <= 1))
    x[index5] = 5 - np.sqrt(4 - 4 * u[index5])
    return x


def split(x, n, shuffle=True):
    n_data = len(x)
    n_data_split = n_data // n
    n_abandoned = n_data % n
    if n_abandoned != 0:
        print(f'Warning: {n_abandoned} samples are abandoned')
    if shuffle:
        x_split = np.random.permutation(x)
    else:
```

4

```python
33          x_split = x.copy()
34          x_split = [x_split[i:i+n_data_split] for i in range(0, n_data,
        n_data_split)]
35          return x_split
36
37
38  def make_train_data(x, n_split, i):
39          x_valid = x[i]
40          x_train = []
41          for j in range(n_split):
42              if j != i:
43                  x_train.extend(x[j])
44          x_train = np.array(x_train)
45          return x_train, x_valid
46
47
48  def gauss_kernel(x, d):
49          k = (1/(2*np.pi)**(d/2)) * np.exp(-(1/2)*x**2)
50          #print((1/(2*np.pi)**(d/2)), k, -(1/2)*x.T.dot(x))
51          return k
52
53
54  def kernel_density(x, h, x_axis, kernel):
55          n = x.shape[0]
56          if len(x.shape) == 1:
57              d = 1
58          else:
59              d = x.shape[1]
60          prob = np.zeros(len(x_axis))
61          for x_i in x:
62              kernel_input = (x_axis - x_i) / h
63              #print(kernel_input)
64              prob += kernel(kernel_input, d)
65          prob = prob / (n * h**d)
66          #print(x)
67          #print(prob)
68          return prob
69
70
71  def estimate_kernel_density(
72          x, n_split, bandwidth_list, kernel,
73          offset=1.0, num=100,
74          path=None,
75          ):
76          x_axis = np.linspace(x.min(), x.max(), num)
77          x_split = split(x, n=n_split, shuffle=True)
78
```

```python
79      n_bandwidth = len(bandwidth_list)
80      n_row = 1
81      n_col = n_bandwidth
82      fig = plt.figure(figsize=(n_col*4, n_row*6))
83
84      lcv_list = []
85      for i, bandwidth in enumerate(bandwidth_list):
86          # calc LCV by likelihood cross validation
87          lcv_list_tmp = []
88          for j in range(n_split):
89              x_train, x_valid = make_train_data(x_split, n_split, j)
90              p = kernel_density(x_valid, bandwidth, x_train, kernel)
91              #lcv = p.sum()
92              lcv = np.log(p).sum()
93              lcv_list_tmp.append(lcv)
94          lcv = np.mean(lcv_list_tmp)
95          lcv_list.append(lcv)
96
97          # plot
98          p = kernel_density(x_train, bandwidth, x_axis, kernel)
99          ax = fig.add_subplot(n_row, n_col, (i+1))
100         ax.set_title(f'$h = {bandwidth}$')
101         ax.hist(x, bins=50, normed=True)
102         ax.plot(x_axis, p)
103         ax.set_ylim([0, 1.0])
104
105     if path:
106         plt.savefig(str(path))
107     plt.show()
108     return lcv_list
109
110
111 def main():
112     # settings
113     n_sample = 3000
114     n_split = 10
115     h_list = [0.01, 0.05, 0.1, 0.5,]  # global
116     #h_list = [0.05, 0.075, 0.1, 0.15]  # local
117     offset = 1.0
118     num = 100
119     fig_path = '../figures/assignment1_result_global.png'
120     np.random.seed(0)
121
122
123     # load data
124     x = generate_data(n_sample)
125     #print('x shape: {}'.format(x.shape))
```

```
126        #plt.hist(x, bins=50)
127        #plt.show()
128
129
130        # train
131        lcv_list = estimate_kernel_density(
132            x, n_split, h_list, gauss_kernel,
133            offset=offset, num=num,
134            path=fig_path,
135            )
136
137        # result
138        form = '{:5.2f}'
139        tab = '  '
140        string_h = '{:3}'.format('h')
141        string_lcv = '{:3}'.format('LCV')
142        for h, lcv in zip(h_list, lcv_list):
143            string_h += tab + form.format(h)
144            string_lcv += tab + form.format(lcv)
145        result = string_h + '\n' + string_lcv
146        print(result)
147
148
149  if __name__ == '__main__':
150        main()
```

**Listings 2:** `assignment2.py`

```python
# -*- coding: utf-8 -*-


import pathlib
import numpy as np
import matplotlib.pyplot as plt


def load_data(n_label=None, n_train=None, n_test=None):
    data_dir = '../data/'
    data_dir = pathlib.Path(data_dir)
    categories = list(range(10))
    train_X = []
    train_y = []
    test_X = []
    test_y = []
    for category in categories[:n_label]:
        # train data
        data_path = data_dir / 'digit_train{}.csv'.format(category)
        data = np.loadtxt(str(data_path), delimiter=',')[:n_train]
        n_data = len(data)
        train_X.extend(data)
        train_y.extend(np.ones(n_data) * category)

        # test data
        data_path = data_dir / 'digit_test{}.csv'.format(category)
        data = np.loadtxt(str(data_path), delimiter=',')[:n_test]
        n_data = len(data)
        test_X.extend(data)
        test_y.extend(np.ones(n_data) * category)
    train_X = np.array(train_X)
    train_y = np.array(train_y)
    test_X = np.array(test_X)
    test_y = np.array(test_y)
    labels = categories[:n_label]
    return train_X, train_y, test_X, test_y, labels


def shuffle(data_X, data_y):
    n_data = len(data_y)
    indices = np.arange(n_data)
    np.random.shuffle(indices)
    data_X_shuffled = data_X[indices]
    data_y_shuffled = data_y[indices]
    return data_X_shuffled, data_y_shuffled

```

```python
47
48   def split(data_X, data_y, n):
49       n_data = len(data_y)
50       n_data_split = n_data // n
51       n_abandoned = n_data % n
52       if n_abandoned != 0:
53           print(f'Warning: {n_abandoned} samples are abandoned')
54       data_X_split = [data_X[i:i+n_data_split] for i in range(0, n_data,
         n_data_split)]
55       data_y_split = [data_y[i:i+n_data_split] for i in range(0, n_data,
         n_data_split)]
56       return data_X_split, data_y_split
57
58
59   def make_train_data(train_X, train_y, n_split, i):
60       train_X_valid = train_X[i]
61       train_y_valid = train_y[i]
62       train_X_train = []
63       train_y_train = []
64       for j in range(n_split):
65           if j != i:
66               train_X_train.extend(train_X[j])
67               train_y_train.extend(train_y[j])
68       train_X_train = np.array(train_X_train)
69       train_y_train = np.array(train_y_train)
70       return train_X_train, train_y_train, train_X_valid, train_y_valid
71
72
73   def knn(train_X, train_y, test_X, k_list, save_memory=False):
74       if save_memory:
75           n_train = train_X.shape[0]
76           n_test = test_X.shape[0]
77           dist_matrix = np.zeros((n_test, n_train))
78           for i in range(n_test):
79               test_X_i = test_X[i]
80               dist_matrix[i, :] = np.sum((train_X - test_X_i[np.newaxis,
         :])**2, axis=1)
81       else:
82           dist_matrix = np.sqrt(
83               np.sum((train_X[None] - test_X[:, None])**2, axis=2)
84           )
85
86       sorted_index_matrix = np.argsort(dist_matrix, axis=1)
87       ret_matrix = None
88       for k in k_list:
89           knn_label = train_y[sorted_index_matrix[:, :k]]
90           label_sum_matrix = None
```

```python
91          for i in range(10):
92              predict = np.sum(np.where(knn_label == i, 1, 0), axis=1)[:, None]
93              if label_sum_matrix is None:
94                  label_sum_matrix = predict
95              else:
96                  label_sum_matrix = np.concatenate(
97                      [label_sum_matrix, predict],
98                      axis=1)
99          if ret_matrix is None:
100             ret_matrix = np.argmax(label_sum_matrix, axis=1)[:, None]
101         else:
102             ret_matrix = np.concatenate([
103                 ret_matrix,
104                 np.argmax(label_sum_matrix, axis=1)[:, None]
105                 ], axis=1)
106     #asert ret_matrix.shape == (len(test_x), len(k_list))
107     return ret_matrix


110 def train(train_X, train_y, k_list, save_memory=False):
111     n_split = len(train_y)
112     n_corrects_list = []
113     for i in range(n_split):
114         train_X_train, train_y_train, train_X_valid, train_y_valid =
        make_train_data(
115             train_X, train_y, n_split, i
116             )
117         y_preds = knn(
118             train_X_train, train_y_train, train_X_valid,
119             k_list, save_memory=save_memory
120             )
121         result = (y_preds == train_y_valid[:, np.newaxis])
122         n_corrects = result.astype(int).sum(axis=0)
123         n_corrects_list.append(n_corrects)
124     n_corrects_list = np.array(n_corrects_list)
125     n_corrects = n_corrects_list.mean(axis=0)
126     return n_corrects


129 def test(train_X, train_y, test_X, test_y, k, labels):
130     n_label = len(labels)
131     confusion_matrix = np.zeros((n_label, n_label), dtype=int)
132     n_data_all = len(test_y)
133     result = {}
134     print('Test')
135
136     preds_all = knn(train_X, train_y, test_X, [k]).reshape(n_data_all)
```

```python
137    #result = (preds_all == test_y)
138    #n_corrects = result.sum(axis=0)
139
140    for label in labels:
141        print(f'Label: {label}\t', end='')
142
143        indices = np.where(test_y == label)[-1]
144        n_data = len(indices)
145        preds = preds_all[indices]
146
147        # make confusion matrix
148        for i in labels:
149            n = (preds == i).sum()
150            confusion_matrix[label, i] = n
151
152        # calc accuracy
153        n_correct = confusion_matrix[label, label]
154        acc = n_correct / n_data
155        print(f'#Data: {n_data}\t#Correct: {n_correct}\tAcc: {acc:.3f}')
156
157        result[label] = {
158            'data': n_data,
159            'correct': n_correct,
160            'accuracy': acc,
161            }
162    result['confusion_matrix'] = confusion_matrix
163
164    # overall score
165    n_crr_all = np.diag(confusion_matrix).sum()
166    acc_all = n_crr_all / n_data_all
167    result['all'] = {
168        'data': n_data_all,
169        'correct': n_crr_all,
170        'accuracy': acc_all,
171        }
172    print(f'All\t#Data: {n_data_all}\t#Correct: {n_crr_all}\tAcc:
       {acc_all:.3f}')
173    print()
174    print('Confusion Matrix:\n', confusion_matrix)
175    print()
176    return result
177
178
179 def print_result_in_TeX_tabular_format(result):
180    labels = list(range(10))
181    print('Scores')
182    for label in labels:
```

11

```python
183         print('{} & {} & {} & {:.3f} \\\\\'.format(
184             label,
185             int(result[label]['data']),
186             int(result[label]['correct']),
187             result[label]['accuracy']
188             ))
189     print()
190     print('Confusion Matrix')
191     for i in labels:
192         print('{}    '.format(i), end='')
193         for j in labels:
194             print(' & {}'.format(int(result['confusion_matrix'][i, j])),
     end='')
195         print(' \\\\')
196     return
197
198
199 def main():
200     # settings
201     k_list = list(range(1, 11, 1))
202     np.random.seed(0)
203     print('Settings')
204     print(f'k Candidates: {k_list}\n')
205
206     # load data
207     train_X, train_y, test_X, test_y, labels = load_data(
208         n_label=None, n_train=None, n_test=None,
209         )
210     _train_X, _train_y = shuffle(train_X, train_y)
211     _train_X, _train_y = split(_train_X, _train_y, n=10)
212
213     # train
214     print('Train')
215     n_corrects = train(_train_X, _train_y, k_list, save_memory=True)
216     print(f'#Correct: {n_corrects}')
217     k_best = np.argmax(n_corrects) + 1
218     print(f'Best k: {k_best}\n')
219
220     # test
221     result = test(train_X, train_y, test_X, test_y, k_best, labels)
222     print_result_in_TeX_tabular_format(result)
223
224
225 if __name__ == '__main__':
226     main()
```