

力学系による将来のパラメータ推定

前多啓一

2019 年 1 月 22 日

1 基本的な方針

注意 1.0.1 (方針). 以下の方針で予測を行う. 与えられたデータは, n 個の観測ポイントにおける関数 $x: \mathbb{R} \rightarrow \mathbb{R}^n$ $t \mapsto (x_1, \dots, x_n)$ の等間隔 τ の時間 t_1, \dots, t_m でのデータである. 推定するのは, k 番目の特定の変数 x_k の将来での動きである.

- $(1, 2, \dots, n)$ のなかから, L 個の数が入っているタプルを s 個選ぶ.
- l 番目のタプルから, 次の値を最小化する $\psi_l: \mathbb{R}^L \rightarrow \mathbb{R}$ を推定する. (ガウス過程回帰)

$$\sum_{i=1}^{m-1} |x_k(t_{i+1}) - \psi_l(x_{l_1}(t_i), x_{l_2}(t_i), \dots, x_{l_L}(t_i))|$$

- 各 ψ_l より 1 ステップの推定 $\tilde{x}_k^l(t + \tau) = \psi_l^l(x_{l_1}(t), \dots, x_{l_L}(t))$ を計算する.
- 集めてできた推定の集合から, カーネル密度推定を行うことで, 確率密度関数 $p(x)$ を推定する.
- 確率密度関数の歪度 γ を計算し, γ が 0.5 以下であれば採用し, $\tilde{x}_k(t + \tau) = \int x p(x) dx$ を推定として確定する. そうでなければ, 以下のように推定値を修正する. 交差検証によりインサンプルエラー δ_l を計算し, それに従って r 個のベストなサンプルを選び出す.

$$\tilde{x}_k(t + \tau) = \sum_{i=1}^r \omega_i \tilde{x}_k^{l_i}(t + \tau)$$

ここで, $\omega_i = \frac{\exp(-\delta_i/\delta_1)}{\sum_j \exp(-\delta_j/\delta_1)}$ である.

定義 1.0.2 (カーネル密度推定). x_1, \dots, x_n を確率密度関数 f をもつ独立同分布からの標本とする. カーネル関数 K , バンド幅 h のカーネル密度推定量とは,

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

基本的に, $K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ を使う. また, 最適なバンド幅として, 以下の値がある.

$$h^* = \frac{c_1^{-2/5} c_2^{1/5} c_3^{-1/5}}{n^{1/5}},$$

where $c_1 = \int x^2 K(x) dx$, $c_2 = \int K(x)^2 dx$, $c_3 = \int (f''(x))^2 dx$.

これについてはカーネル密度推定が `scipy` に標準搭載されているのでそちらを援用.

2 ガウス過程回帰について

Bishop を参照 [2] しながら, ガウス過程回帰について復習する.

- 推定の仮定

$y = \mathbf{w}^T \phi(\mathbf{x})$ とし, パラメータ \mathbf{w} がガウス分布に従うと仮定する.

すなわち, 任意の $\mathbf{x}_1, \dots, \mathbf{x}_n$ に対し, $\mathbf{y} = \Phi \mathbf{w}$ はガウス分布に従う. このことから, \mathbf{y} は無限次元のガウス分布に従う, などとも言われる. ただし, $\Phi = (\phi(\mathbf{x}_i))_{i=1, \dots, n}$ は計画行列である.

- 与えられるデータ (サンプル)

$\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^n$ および $t_1, \dots, t_n \in \mathbb{R}$ ただし, $t_n = y_n + \varepsilon_n$ であるとする. ε_n はノイズで, ガウス分布に従うとする.

- 推定するもの

新しい入力 \mathbf{x}_{n+1} が与えられたときの出力 t_{n+1} の確率分布を推定する. すなわち,

$$p(t_{n+1} | \mathbf{x}_{n+1}, \mathbf{x}_1, \dots, \mathbf{x}_n, t_1, \dots, t_n) = N(t_{n+1} | m, \sigma^2)$$

における m と σ^2 を推定する.

定理 2.0.1. 以下のようにカーネル関数のグラム行列を定義する.

$$K = (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j}$$

さらに, 以下のように置く.

$$\mathbf{t} = \begin{pmatrix} t_1 \\ \vdots \\ t_n \end{pmatrix}, \quad \mathbf{k} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}) \\ \vdots \\ k(\mathbf{x}_n, \mathbf{x}) \end{pmatrix}$$

最適な推定は, 以下の通り.

$$m = \mathbf{k}^T (K + \sigma_n^2 I)^{-1} \mathbf{t}$$
$$\sigma^2 = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}^T (K + \sigma_n^2 I)^{-1} \mathbf{k}$$

3 コード

以上を踏まえ, 以下のようにコードを組んだ (参考 [1])

ソースコード 1 main.py

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 plt.style.use('ggplot')
5 from scipy.stats import gaussian_kde
6
7 # 得られている情報の数
8 initial = 28
```

```

9 # 予測に使う数
10 n = 10
11 # 予測に使うデータの個数
12 m = 4
13 # 予測の回数
14 p = 30
15 # 予測に使うデータの深さ
16 q = 10
17 # カーネル関数の情報
18 a = 0.05
19 b = 0.05
20 beta = 100
21
22
23 # 予測する株の名前
24 target = '1812_JT_Equity'
25
26 similar = ['1801_JT_Equity', '1802_JT_Equity', '1803_JT_Equity', '1812_JT_Equity', '1820_JT_
Equity', '1821_JT_Equity', '1824_JT_Equity', '1833_JT_Equity', '1860_JT_Equity', '1893_JT_
Equity']
27 # 推定のために 1/1000倍
28 df = pd.read_csv('./Data/Open.csv', usecols = similar)/1000
29
30 df_real = df[0:initial+p]
31 df = df[0:initial]
32
33 # リストを, ある値との距離順に並べる関数
34 def pointsort(arr, val):
35     return [y[1] for y in sorted([(abs(x-val), x) for x in arr])]
36
37 # random に n 個の配列を作る
38 def pickup(arr, n):
39     arr2 = arr[:]
40     result = []
41     for i in range(n):
42         x = arr2[int(len(arr2) * np.random.rand())]
43         result.append(x)
44         arr2.remove(x)
45     return result
46
47 # a を含まないランダムな配列を n 個作る
48 def pickup2(arr, a, n):
49     arr2 = arr[:]
50     arr2.remove(a)
51     return pickup(arr2, n-1)
52
53 # 多変数のグラム行列を作る関数
54 def gram(f, x):
55     result = []
56     for i in x:
57         smallresult = []
58         for j in x:
59             smallresult.append(f(i, j))
60         result.append(smallresult)
61     return np.array(result)
62
63 # ベクトルを作る関数

```

```

64 def k(f,new,x):
65     result = []
66     for i in x:
67         result.append(f(new,i))
68     return np.array(result)
69
70 # カーネル推定のクラス
71 class GaussianKernel(object):
72     def __init__(self, params):
73         assert np.shape(params) == (2,)
74         self.__params = params
75
76     def __call__(self, x1, x2):
77         return self.__params[0] * np.exp(-0.5 * self.__params[1] * np.linalg.norm(x1 - x2)
78             ** 2)
79
80     # 以下の関数はparameter の推定のための用意
81     def get_params(self):
82         return np.copy(self.__params)
83
84     def derivatives(self, x1, x2):
85         delta_1 = -0.5 * sq_diff * delta_0 * self.__params[0]
86         return (delta_0, delta_1)
87
88     def delta0(self,x1,x2):
89         sq_diff = np.linalg.norm(x1 - x2) ** 2
90         return np.exp(-0.5 * self.__params[1] * sq_diff)
91
92     def delta1(self,x1,x2):
93         sq_diff = np.linalg.norm(x1 - x2) ** 2
94         return -0.5 * sq_diff * self.delta0(x1,x2) * self.__params[0]
95
96     def update_parameters(self, updates):
97         assert np.shape(updates) == (2,)
98         self.__params += updates
99
100 # GPR 用のクラス
101 class GaussianProcessRegression(object):
102     def __init__(self, kernel, beta=1.):
103         self.kernel = kernel
104         self.beta = beta
105
106     def fit(self, x, y):
107         self.x = x
108         self.y = y
109         Gram = gram(self.kernel,x)
110         self.covariance = Gram + np.identity(len(x)) / self.beta
111         self.precision = np.linalg.inv(self.covariance)
112
113     # Kernel function の parameter 推定 / 今回は綺麗に収束しないので, 要検討
114     def fit_kernel(self, x, y, learning_rate=0.1, iter_max=100):
115         for i in range(iter_max):
116             params = self.kernel.get_params()
117             self.fit(x, y)
118             grad0 = gram(self.kernel.delta0,x)
119             grad1 = gram(self.kernel.delta1,x)
120             gradients = [grad0,grad1]

```

```

120         updates = np.array(
121             [-np.trace(self.precision.dot(grad)) + y.dot(self.precision.dot(grad).dot(
122                 self.precision).dot(y)) for grad in gradients])
122         print(updates)
123         self.kernel.update_parameters(learning_rate * updates)
124         if np.allclose(params, self.kernel.get_params()):
125             break
126         else:
127             print("parameters may not have converged")
128
129     def predict_dist(self, new):
130         K = k(self.kernel, new, self.x)
131         mean = K.dot(self.precision).dot(self.y)
132         var = self.kernel(new, new) + 1 / self.beta - np.sum(K.dot(self.precision) * K)
133         return mean.ravel(), np.sqrt(var.ravel())
134
135     # 期待値を算出する関数
136     def expectation(x,y):
137         y = y / sum(y)
138         return sum(x*y)
139
140     # kernel density estimation をしたのち、期待値を算出する
141     def kde_process(data_list):
142         kde_model = gaussian_kde(data_list)
143         y = kde_model(data_list)
144         # x_grid = np.linspace(min(data_list), max(data_list), num=100)
145         # estimy = kde_model(x_grid)
146         # データを正規化したヒストグラムを表示する用
147         # weights = np.ones_like(data_list)/float(len(data_list))
148         # print("mean:", pandas.mean())
149         # plt.figure(figsize=(14,7))
150         # plt.plot(x_grid, estimy)
151         # plt.hist(data_list, alpha=0.3, bins=20, weights=weights)
152         # plt.show()
153         skew = pd.Series(y).skew()
154         if abs(skew) < 0.1:
155             return expectation(data_list,y)
156         else:
157             data_list2=pointsort(data_list,np.average(data_list))[0:int(len(data_list)/3)]
158             return expectation(data_list2,kde_model(data_list2))
159
160     # データフレームdfの中から、target の株の mean と sd を推定する関数
161     def onetimeestimation(df,target,n,m,a,b,beta):
162         # target の情報 (1 つだけずらして取得)
163         y = df[target].values[1:]
164         y = y[:n-1]
165         kernel = GaussianKernel(params=np.array([a, b]))
166         result = np.array([])
167         result_sd = np.array([])
168         for i in range(n):
169             x = df[pickup2(similar, target, m)].values
170             # 推定前のx(matrix)
171             x = x[:n]
172             x,x_test = x[:-1],x[-1]
173             regression = GaussianProcessRegression(kernel=kernel, beta=beta)
174             # regression.fit_kernel(x, y, learning_rate=0.1, iter_max=10000)
175             regression.fit(x,y)

```

```

176     pred_y, pred_y_sd = regression.predict_dist(x_test)
177     result = np.append(result, pred_y)
178     result_sd = np.append(result_sd, pred_y_sd)
179     mean = kde_process(result)
180     sd = np.average(result_sd)
181     return mean, sd
182
183 # 同業種リストsimilar のすべての1ステップを推定する関数
184 def onetimeallestimation(df, similar, n, m, q, a, b, beta):
185     result = []
186     result_sd = []
187     df = df[-q:]
188     for terget in similar:
189         result.append(onetimeestimation(df, terget, n, m, a, b, beta)[0])
190         result_sd.append(onetimeestimation(df, terget, n, m, a, b, beta)[1])
191     return result, result_sd
192
193 # 複数回の推定を行う
194 def manytimeestimation(df, similar, n, m, p, q, a, b, beta):
195     error = df[0:0].copy()
196     for i in range(p):
197         mean, sd = onetimeallestimation(df, similar, n, m, q, a, b, beta)
198         mean = pd.Series(mean, index=similar, name='pred'+str(i))
199         sd = pd.Series(sd, index=similar, name=i)
200         df = df.append(mean)
201         error = error.append(sd)
202     return df, error
203
204 # グラフの出力のための関数
205 def makegraph(estimation, sd, real):
206     x_grid = np.array(range(len(estimation)))
207     plt.figure(figsize=(14, 7))
208     plt.errorbar(x_grid, estimation, sd, fmt='ro-')
209     plt.plot(x_grid, real, 'go-')
210     plt.show()
211
212 result, sd = manytimeestimation(df, similar, n, m, p, q, a, b, beta)
213 estimation = result[terget].values[initial:initial+p]
214 sd = sd[terget].values[0:p]
215 real = df_real[terget].values[initial:initial+p]
216
217 # print(onetimeestimation(df, terget, n, m, a, b, beta))
218 # print(onetimeallestimation(df, similar, 10, 3, 10))
219 makegraph(estimation, sd, real)

```

4 結果

以下の同業種の株で推定を行った。

推定グループ (similar) : 業種グループ:Engineering 業種サブグループ:Building&Construc-Misc

'1801 JT Equity', '1802 JT Equity', '1803 JT Equity', '1812 JT Equity', '1820 JT Equity',
 '1821 JT Equity', '1824 JT Equity', '1833 JT Equity', '1860 JT Equity', '1893 JT Equity'

推定した株 (target) : '1812 JT Equity'

推定に費やした回数 (n) : 10 回

推定に使った株の数 (m) : 3 ずつ

1 回のカーネル密度推定は以下のようになる.

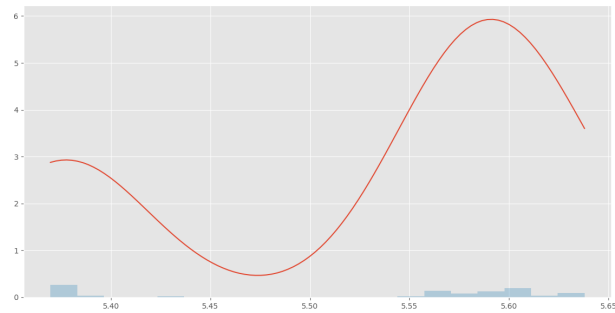


図 1 カーネル密度推定

以下は, 10 回すべての株を推定し, 将来予測を行なった結果である.

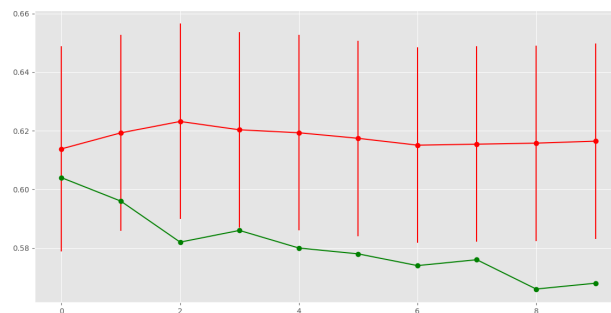


図 2 10 回予測

参考文献

- [1] Qiita PRML 第 6 章 ガウス過程による回帰 Python 実装 <https://qiita.com/ctgk/items/4c4607edf15072cddc46>
- [2] Christopher M. Bishop “Pattern Recognition and Machine Learning” 2013