

# 力学系による将来のパラメータ推定

前多啓一

2019 年 3 月 18 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  $x: \mathbb{R} \rightarrow \mathbb{R}^n \quad t \mapsto (x_1, \dots, x_n)$  の等間隔  $\tau$  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  $\psi_l$  より 1 ステップの推定  $\tilde{x}_k^l(t + \tau) = \psi_k^l(x_{l_1}(t), \dots, x_{l_L}(t))$  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  $\gamma$  を計算し,  $\gamma$  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  $\delta_l$  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.

**定義 2.0.1** (カーネル関数 (正定値カーネル)).  $\Omega$  を集合とし,  $k: \Omega \times \Omega \rightarrow \mathbb{R}$  を写像とする.  $k$  が  $\Omega$  の正定値カーネルであるとは, 次の 2 つを満たすことを言う.

- (1)  $k$  は対称. すなわち,  $k(x, y) = k(y, x)$  である.
- (2)  $k$  の  $n$  次元グラム行列  $(k(x_i, x_j))_{i,j}$  が半正定値である.

- 推定の仮定

関数  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^n$  に対し,  $y = \mathbf{w}^T \phi(\mathbf{x})$  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.  $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}: (x, x') \mapsto k(x, x') = \phi(x)^T \phi(x')$  はカーネル関数である.

- 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.  $\varepsilon_n$  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  $\sigma^2$  を推定する.

**定理 2.0.2.** 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.

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

## 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]). 比較のため, 線形回帰および LASSO 回帰による計算も行う.

---

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu March 7 10:04:54 2019
4
5  @author: maeta
6  """
7
8  import pandas as pd
9  import numpy as np
10 import matplotlib.pyplot as plt
11 plt.style.use('ggplot')
12 from scipy.stats import gaussian_kde
13 from sklearn.gaussian_process import kernels as sk_kern
14 from sklearn.gaussian_process import GaussianProcessRegressor
15 from sklearn import linear_model
16
17 #予測期間
18 start = 20160107
19 end = 20161230
20
21 similar = ['1801_JT_Equity',
22 '1802_JT_Equity',
23 '1803_JT_Equity',
24 '1812_JT_Equity',
25 '1820_JT_Equity',
26 '1821_JT_Equity',
27 '1824_JT_Equity',
28 '1833_JT_Equity',
29 '1860_JT_Equity',
30 '1893_JT_Equity']
31
32 # 推定のために 1日リターンに
33 df_open = pd.read_csv('Open.csv',usecols =similar)
34 df_close = pd.read_csv('Close.csv',usecols =similar)
35 df = (df_close - df_open) /df_open
36
37 date = pd.read_csv('Date.csv')
38 po_start = int(date[date['Date']==start]['No'])
39 po_end = int(date[date['Date']==end]['No'])
40 index = date['Index'][po_start:po_end+1]
41
42 ##GPR による方法
43 # df の similar に入っている株を, p 個ペアで q 回ずつ, n の深さで予想する.
44 class GPReestimation:
45     def __init__(self, df, similar, n, p, q):
46         self.df = df
47         self.similar = similar
48         self.n = n
49
50         self.p = p
51         self.q = q
52
53     # リストを, ある値との距離順に並べる関数
54     def pointsort(self,arr,val):
55         return [y[1] for y in sorted([(abs(x-val),x) for x in arr])]
56

```

```

57     # random に n 個の配列を作る
58     def pickup(self,arr,n):
59         arr2 = arr[:]
60         result = []
61         for i in range(n):
62             x = arr2[int(len(arr2) * np.random.rand())]
63             result.append(x)
64             arr2.remove(x)
65         return result
66
67     # a を含まないランダムな配列を n 個作る
68     def pickup2(self,arr,a,n):
69         arr2 = arr[:]
70         arr2.remove(a)
71         return self.pickup(arr2,n-1)
72
73     # GPR
74     def GPR_fit(self,x_train,y_train,x_test):
75         kernel = sk_kern.RBF(1.0, (1e-3, 1e3)) + sk_kern.ConstantKernel(1.0, (1e-3, 1e3)) +
76             sk_kern.WhiteKernel()
77         clf = GaussianProcessRegressor(
78             kernel=kernel,
79             alpha=1e-10,
80             optimizer="fmin_l_bfgs_b",
81             n_restarts_optimizer=20,
82             normalize_y=True)
83         clf.fit(x_train,y_train)
84         pred_mean, pred_std = clf.predict(x_test, return_std=True)
85         return pred_mean,pred_std
86
87     # 期待値を算出する関数
88     def expectation(self,x,y):
89         y = y / sum(y)
90         return sum(x*y)
91
92     # kernel density estimation をしたのち、期待値を算出する
93     def kde_process(self,data_list):
94         kde_model = gaussian_kde(data_list)
95         y = kde_model(data_list)
96         skew = pd.Series(y).skew()
97         if abs(skew) < 0.1:
98             return self.expectation(data_list,y)
99         else:
100             data_list2=self.pointsort(data_list,np.average(data_list))[0:int(len(data_list)
101                 )/3)]
102             return self.expectation(data_list2,kde_model(data_list2))
103
104     def onetimeestimation(self,i,terget):
105         # terget の情報 (1 つだけずらして取得)
106         y_train = self.df[terget].values[i-self.n:i]
107         result = np.array([])
108         result_sd = np.array([])
109         j = 0
110         while j < self.q:
111             try:
112                 x = self.df[self.pickup2(self.similar, terget, self.p)].values
113                 x_train,x_test = x[i-self.n:i],[x[i]]

```

```

112         pred_y ,pred_y_sd = self.GPR_fit(x_train,y_train,x_test)
113         result = np.append(result,pred_y)
114         result_sd = np.append(result_sd,pred_y_sd)
115         j += 1
116     except:
117         j += 1
118     mean = self.kde_process(result)
119     err = (self.df[target].values[i+1]-mean)**2
120     sd = np.average(result_sd)
121     return mean,sd,err
122
123     # 同業種リストsimilar のすべての 1 ステップを推定する関数
124     def onetimeallestimation(self,i):
125         result = []
126         result_sd = []
127         err = 0
128         for target in self.similar:
129             result.append(self.onetimeestimation(i,target)[0])
130             result_sd.append(self.onetimeestimation(i,target)[1])
131             err += self.onetimeestimation(i,target)[2]
132         return result,result_sd,err
133
134     ## 線形回帰による方法
135     class Linearestimation:
136         def __init__(self, df, similar, n):
137             self.df=df
138             self.similar = similar
139             self.n = n
140
141         def linearregression(self,train_x,train_y,test_x):
142             clf = linear_model.LinearRegression()
143             clf.fit(train_x,train_y)
144             return np.dot(clf.coef_,test_x)
145
146         def linearestimation(self,target,n,i):
147             train = self.similar[:]
148             train.remove(target)
149             train_x = self.df[train].values[i-n:i]
150             train_y = self.df[target].values[i-n:i]
151             test_x = self.df[train].values[i:i+1][0]
152             result = self.linearregression(train_x,train_y,test_x)
153             err = (result - self.df[target].values[i+1])**2
154             return result, err
155
156         def linearonetimeallestimation(self,i):
157             result = np.array([])
158             err = 0
159             for target in similar:
160                 result = np.append(result,self.linearestimation(target,self.n,i)[0])
161                 err += self.linearestimation(target,self.n,i)[1]
162             return result,err
163
164     ## LASSO による方法
165     class Lassoestimation:
166         def __init__(self, df, similar, n):
167             self.df = df
168             self.similar = similar

```

```

169         self.n = n
170
171     def lassoregression(self,train_x,train_y,test_x):
172         lasso = linear_model.Lasso(alpha=0.00001)
173         lasso.fit(train_x,train_y)
174         return np.dot(lasso.coef_,test_x)
175
176     def lassoestimation(self,terget,n,i):
177         train = self.similar[:]
178         train.remove(terget)
179         train_x = self.df[train].values[i-n:i]
180         train_y = self.df[terget].values[i-n:i]
181         test_x = self.df[train].values[i:i+1][0]
182         result = self.lassoregression(train_x,train_y,test_x)
183         err = (result - self.df[terget].values[i+1])**2
184         return result, err
185
186     def lassoonetimeallestimation(self,i):
187         result = np.array([])
188         err = 0
189         for terget in similar:
190             result = np.append(result,self.lassoestimation(terget,self.n,i)[0])
191             err += self.lassoestimation(terget,self.n,i)[1]
192         print(err)
193         return result,err
194
195 # GPRcheck
196 PortRet = []
197 GPR = GPReestimation(df,similar,10,3,5)
198 err = 0
199 for i in range(po_start,po_end+1):
200     tmpdf = df[i:i+1]
201     result = GPR.onetimeallestimation(i)[0]
202     result_sd = GPR.onetimeallestimation(i)[1]
203     err += GPR.onetimeallestimation(i)[2]
204     long = similar[np.argmax(result)]
205     short = similar[np.argmin(result)]
206     portreturn = tmpdf[long].values[-1] - tmpdf[short].values[-1]
207     PortRet.append(portreturn)
208     print(str(i)+"...finished")
209     print(err)
210 PortRet = pd.DataFrame(PortRet)
211 PortRet.index = index
212 PortRet.columns = ['Ret-GPR']
213 PortRet.to_csv('Result.csv')
214
215 # linearcheck
216 PortRet = []
217 Line = Linearestimation(df,similar,10)
218 err = 0
219 for i in range(po_start,po_end+1):
220     tmpdf = df[i-11:i]
221     result = Line.linearonetimeallestimation(i)[0]
222     err += Line.linearonetimeallestimation(i)[1]
223     long = similar[np.argmax(result)]
224     short = similar[np.argmin(result)]
225     portreturn = tmpdf[long].values[-1] - tmpdf[short].values[-1]

```

```

226     PortRet.append(portreturn)
227     print(str(i)+"...finished")
228     print(err)
229 print(err)
230 PortRet = pd.DataFrame(PortRet)
231 PortRet.index = index
232 PortRet.columns = ['Ret-Linear']
233 PortRet.to_csv('Result-Linear.csv')
234
235 # Lassocheck
236 PortRet = []
237 Lasso = Lassoestimation(df,similar,10)
238 err = 0
239 for i in range(po_start,po_end+1):
240     tmpdf = df[i-11:i]
241     result = Lasso.lassoonetimeallestimation(i)[0]
242     err += Lasso.lassoonetimeallestimation(i)[1]
243     long = similar[np.argmax(result)]
244     short = similar[np.argmin(result)]
245     portreturn = tmpdf[long].values[-1] - tmpdf[short].values[-1]
246     PortRet.append(portreturn)
247     print(str(i)+"...finished")
248     print(err)
249 PortRet = pd.DataFrame(PortRet)
250 PortRet.index = index
251 PortRet.columns = ['Ret-Lasso']
252 PortRet.to_csv('Result-Lasso.csv')

```

---

## 4 結果

以下の同業種の株で推定を行った.

group1 : 業種グループ:Engineering 業種サブグループ:Building&Construc-Misc

group2 : 業種グループ:BasicMaterials 業種サブグループ:Chemicals

group3 : 業種グループ:Consumer,Noncyclical 業種サブグループ:food

group4 : 業種グループ:Machinery diversified 業種サブグループ:Machinery-generalindust

	Linear	Lasso	GPR
group1	0.75427	0.75644	0.6943
group2	3.2203	2.6910	2.0573
group3	1.7092	0.5125	0.3745
group4	0.6737	0.5925	0.51308

表1 二乗誤差の比較

## 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