

统计信号处理第二次大作业报告

无88 刘子源 2018010895

一、题目描述

给定特征方程

$$\lambda^3 + a\lambda^2 + b\lambda + c = 0$$

且

$$\lambda_1 \geq \lambda_2 \geq \lambda_3$$

令

$$p_i = \frac{\lambda_i}{\lambda_1 + \lambda_2 + \lambda_3}$$

试利用上述特征多项式的系数近似表示下面的表达式

$$H = -\sum_{i=1}^3 p_i \log_3 p_i$$

二、公式化简

将特征多项式进行因式分解：

$$(\lambda - \lambda_1)(\lambda - \lambda_2)(\lambda - \lambda_3) = 0$$

可得到一元三次方程的韦达定理：

$$\lambda_1 + \lambda_2 + \lambda_3 = -a$$

$$\lambda_1 \lambda_2 + \lambda_2 \lambda_3 + \lambda_1 \lambda_3 = b$$

$$\lambda_1 \lambda_2 \lambda_3 = -c$$

由此得到 $a < 0, b > 0, c < 0$

将 $p_i = \frac{\lambda_i}{\lambda_1 + \lambda_2 + \lambda_3}$ 代入原方程，得到关于p的方程：

$$p^3 - p^2 + \frac{b}{a^2}p - \frac{c}{a^3} = 0$$

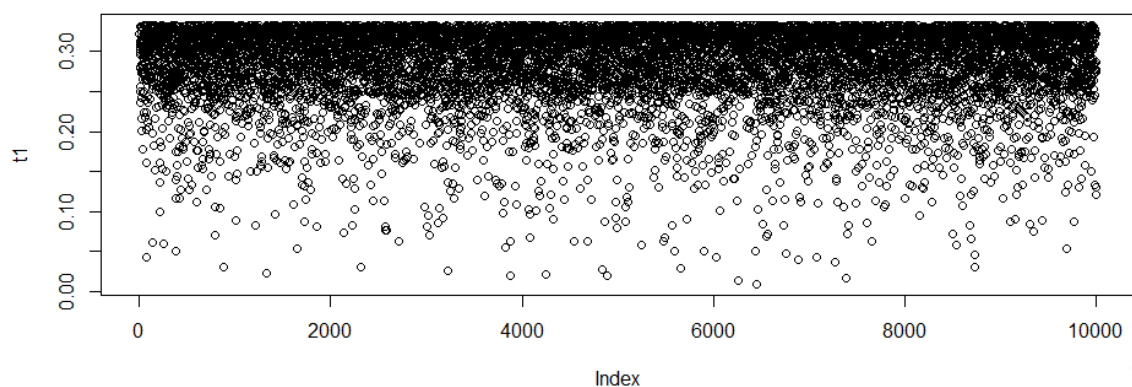
这样可以将a、b、c三个参数化简为 $\frac{b}{a^2}$ 和 $\frac{c}{a^3}$ 两个参数，接下来要寻找H与它的关系。

三、探索性分析

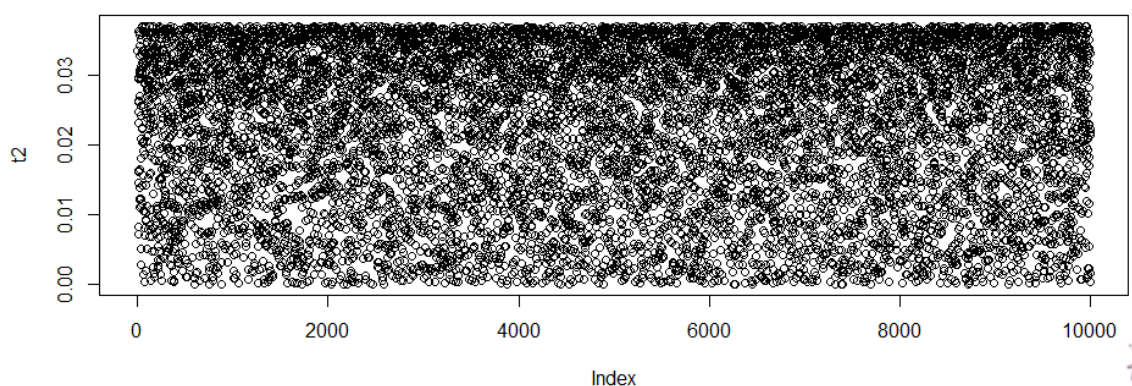
给定参数a, b, c求出 λ 较为困难，因为非常容易出现复数解得情况，观察到 $0 < p_i < 1$ ，其数量级不会受 λ_i 的影响，故可以随机生成 λ ，进而求得参数a, b, c。

首先利用R语言对数据进行探索性分析，随机生成10000个特征方程的解，每个 λ 的取值在0 ~ 10000之间，观察 $t_1 = \frac{b}{a^2}, t_2 = \frac{c}{a^3}, H$ 的分布。

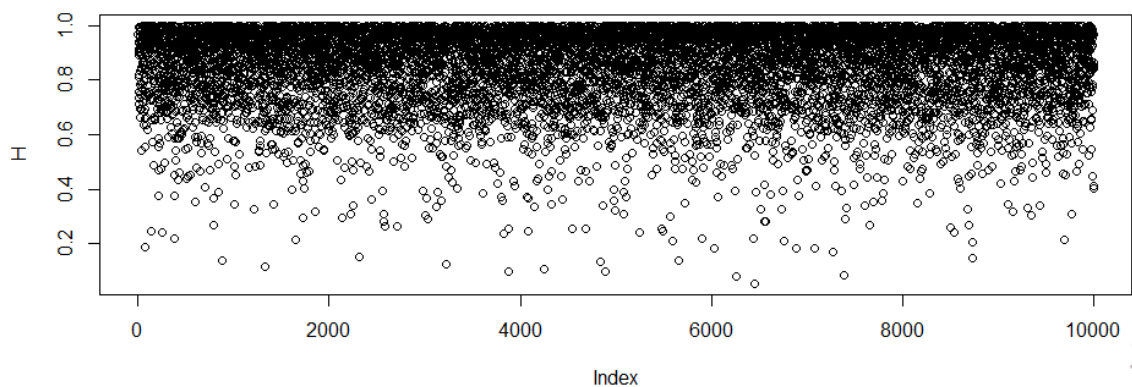
t_1 分布在大约[0, 0.35]的区间上，且在0附近分布稀疏，在0.3附近分布密集：



t_2 分布在大约 $[0, 0.04]$ 的区间上，且分布比较均匀：

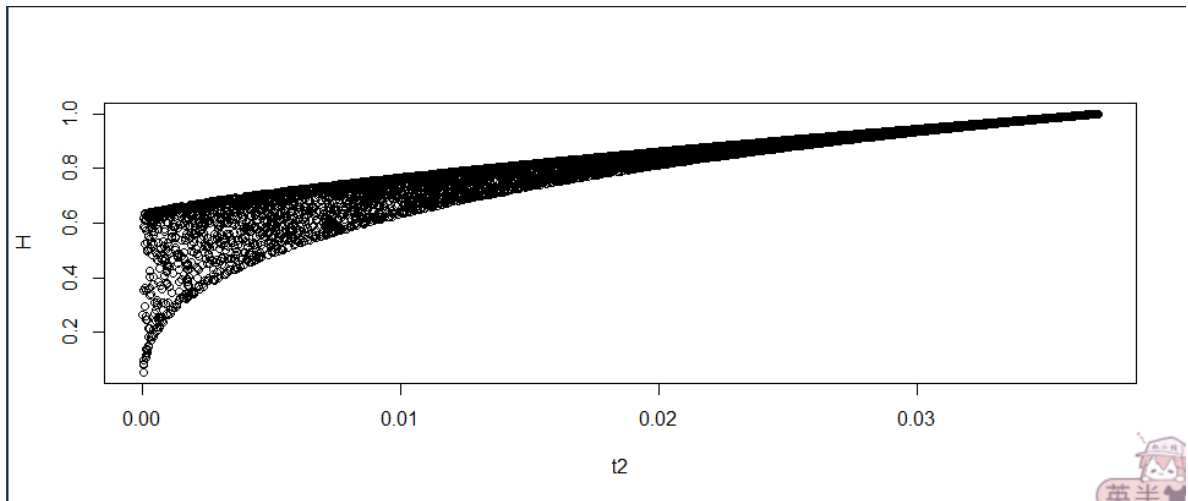
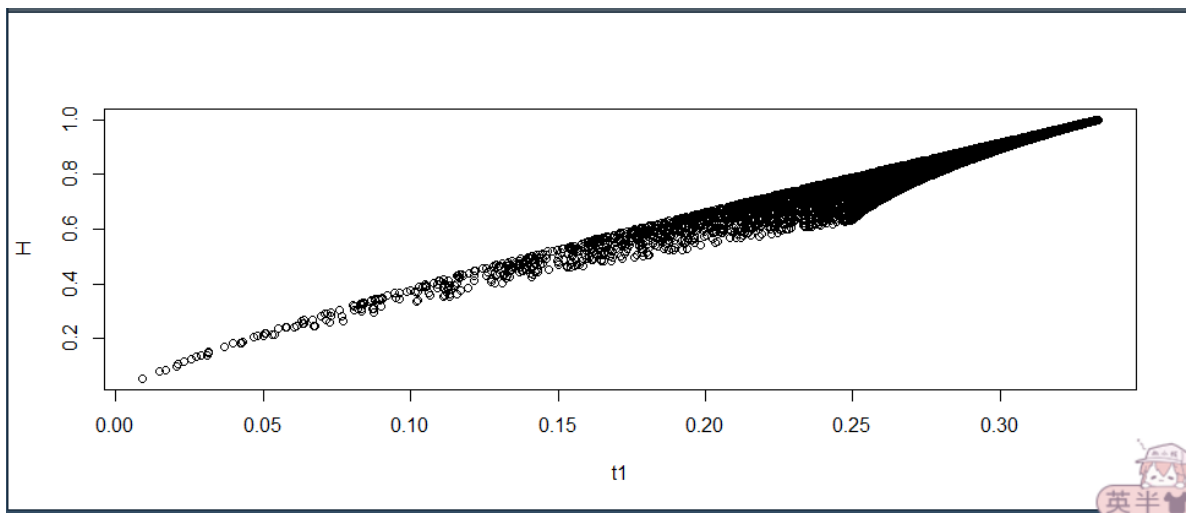


要求的H分布在区间 $[0, 1]$ 之间，在靠近0的地方分布稀疏，在靠近1的地方分布密集：



我们初步发现，H与变量 t_1, t_2 有着相近的数量级，因此在后续拟合过程中无须做方差的归一化处理。

接下来，首先观察H与 t_1, t_2 是否存在着最简单的线性关系，画出H vs t_1 图像和H vs t_2 图像：



可以明显的看到，H与 t_1, t_2 确实存在着线性关系，我们可以以线性回归模型为基础，通过增添高次项、对数项等的方法来完善模型。

四、初步建模及诊断

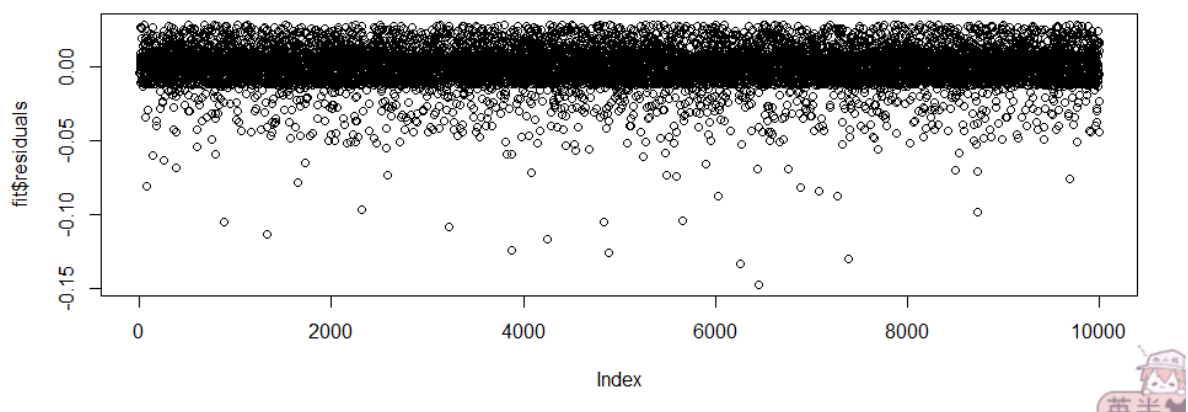
建立线性回归模型

$$H = \beta_0 + \beta_1 t_1 + \beta_2 t_2$$

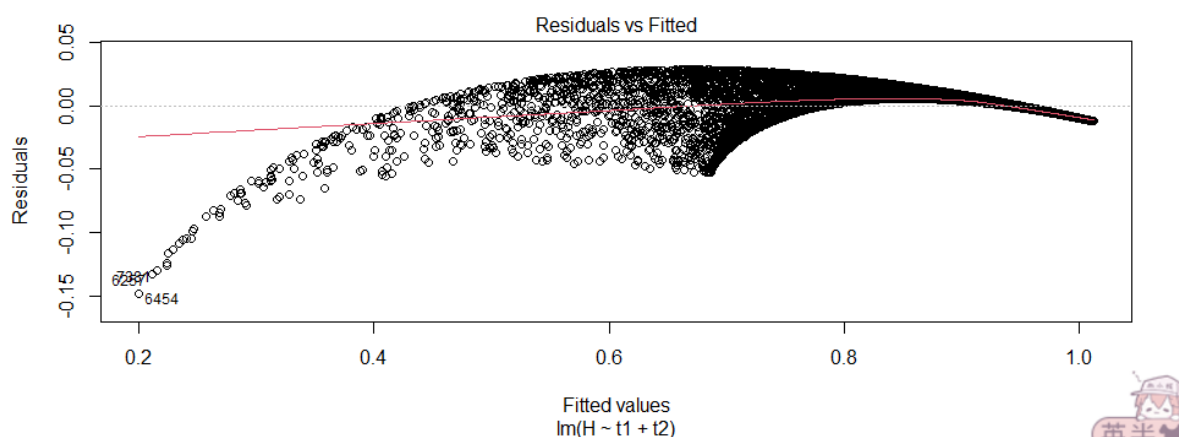
利用最小二乘得到如下估计：

$$\hat{H} = b_0 + b_1 t_1 + b_2 t_2$$

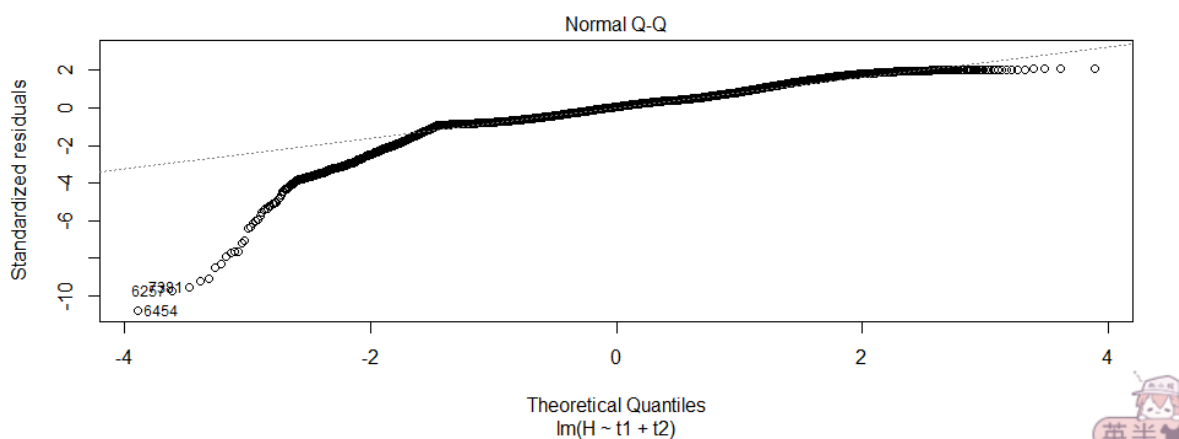
平均误差达到了0.0098，达到了0.01的标准，但是需要进一步检验残差，画出残差图，可以发现存在较多的离群值：



进一步画出residuals vs \hat{H} 图像，可以明显观察到模型存在严重的异方差性，在 H 较小时误差很大，模型的鲁棒性较差。



可以通过正态性检验来分析模型出现异方差性的原因，画出残差的Q-Q图，可以发现模型存在较为严重的左偏和拖尾性，这时可以通过引入对数项来改善。



五、改进模型及诊断

经过多次用R语言对数据进行拟合后，决定模型中使用如下7个解释变量：

$$t_1 = \frac{b}{a^2}, \quad t_2 = \frac{c}{a^3}$$

对数项 $\log(t_1), \log(t_2)$

交叉项 $t_1 t_2, t_1 \log(t_1), t_2 \log(t_2)$

建立回归模型

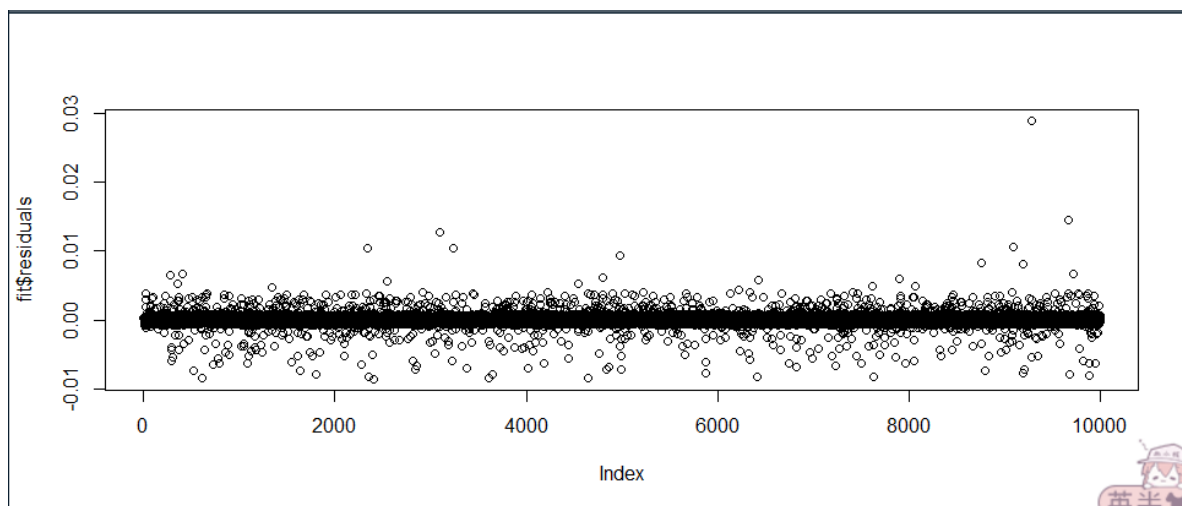
$$H = \beta_0 + \beta_1 t_1 + \beta_2 t_2 + \beta_3 \log(t_1) + \beta_4 \log(t_2) + \beta_5 t_1 \log(t_1) + \beta_6 t_2 \log(t_2) + \beta_7 t_1 t_2$$

利用最小二乘得到如下估计：

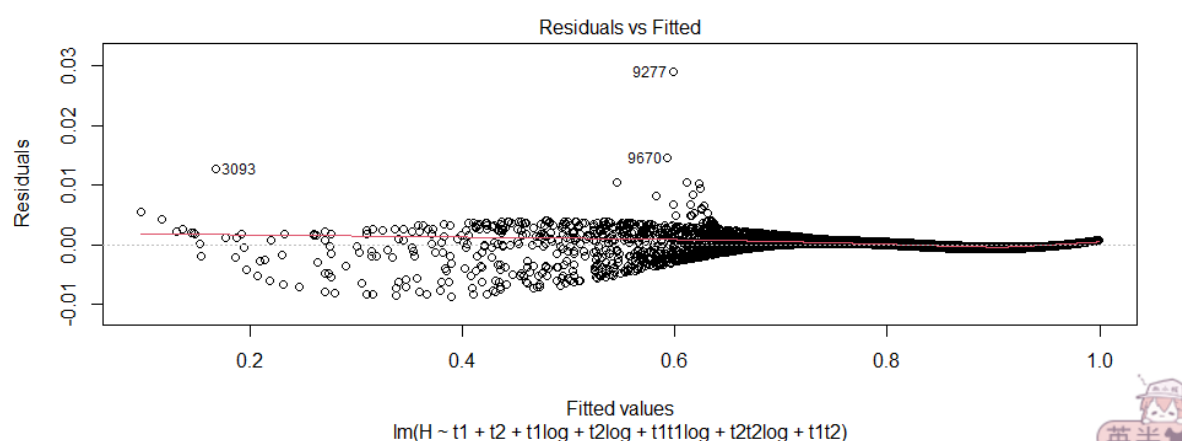
$$\hat{H} = b_0 + b_1 t_1 + b_2 t_2 + b_3 \log(t_1) + b_4 \log(t_2) + b_5 t_1 \log(t_1) + b_6 t_2 \log(t_2) + b_7 t_1 t_2$$

此时平均误差可以降低至 10^{-4} 量级，满足要求。

做出残差图，相比于最初建立的模型，离群现象已大大改善。



画出residuals vs \hat{H} 图像，观察到异方差性已得到很大的改善，在H较小时误差也会稳定在0.01以内，模型的稳定性很好。



所以我选择此模型为最终模型。

六、测试结果

用matlab对模型进行测试验证。

step1

首先随机生成50000个数据对模型参数 $\beta = [\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7]$ 进行估计，得到 β 的估计值

$$\hat{\beta} = [0.0701, 1.3264, 0.3147, 0.0034, 0.0047, -0.7773, -2.0438, -3.2013]$$

即H的估计为：

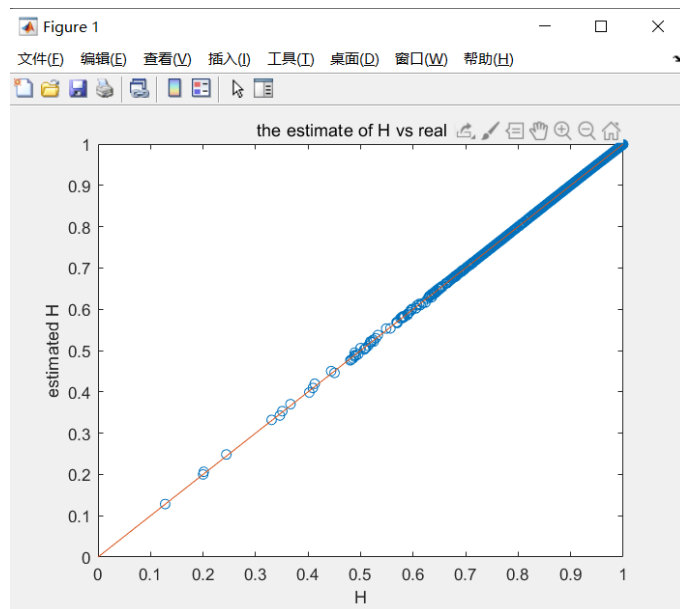
$$\hat{H} = 0.0701 + 1.3264t_1 + 0.3147t_2 + 0.0034\log(t_1) + 0.0047\log(t_2) - 0.7773t_1\log(t_1) - 2.0438t_2\log(t_2) - 3.2013t_1t_2$$

step2

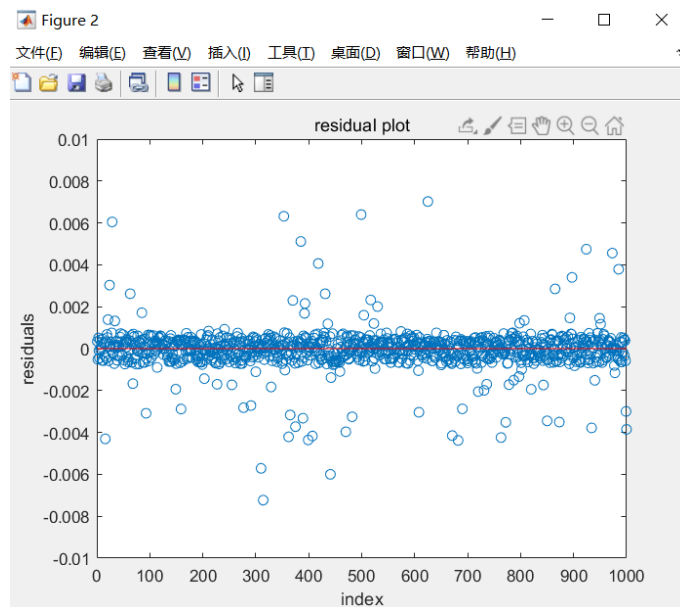
随机生成1000个测试点，得到平均误差 $\overline{error} = 5.496 \times 10^{-4}$

```
>> error
error =
5.4958e-04
```

做出H的估计值 vs H真实值图像，点几乎都分布在直线 $y = x$ 附近。



画出残差图，残差已很好地近似为稳定的白噪声。



step 3

作业要求，输入系数 a, b, c ，输出 H 的估计量，但由于 a, b, c 有较多的限制，容易出现复数解的情况，所以改为输入特征方程的解 $\lambda_1, \lambda_2, \lambda_3$ 的解，由方程的解倒推出参数 a, b, c ，再由参数 a, b, c 得到 H 的估计量 \hat{H} ，并与 H 的真实值进行比较：

七、复杂度分析

$$\begin{aligned} \text{b}[19] & \Rightarrow \text{Simplify(H)} \\ \text{Out}[19] & = -\left(\left(-2 \left(2^{2/3} a^3 - 6 \cdot 2^{1/3} b - 2a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} + \left(-4a^3 + 18ab - 54c + 6\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{2/3} \right) \right. \\ & \quad \left. \log \left[-2 \cdot 2^{1/3} a^3 + 6 \cdot 2^{1/3} b + 2a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} - \left(-4a^3 + 18ab - 54c + 6\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{2/3} \right] \right. \\ & \quad \left. \left(6a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} \right) + \left(2 \cdot 2^{1/3} (1+i\sqrt{3}) a^2 + 2^{1/3} (-6-6i\sqrt{3}) b + \right. \right. \\ & \quad \left. \left. 4a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} + (1-i\sqrt{3}) \left(-4a^3 + 18ab - 54c + 6\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{2/3} \right) \right) \\ & \quad \log \left[\left(2 \cdot 2^{1/3} (1+i\sqrt{3}) a^2 + 2^{1/3} (-6-6i\sqrt{3}) b + 4a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} + \right. \right. \\ & \quad \left. \left. (1-i\sqrt{3}) \left(-4a^3 + 18ab - 54c + 6\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{2/3} \right) \right] \left/ \left(12a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} \right) + \right. \\ & \quad \left. \left(2 \cdot 2^{1/3} (1-i\sqrt{3}) a^2 + 6i \cdot 2^{1/3} (1+\sqrt{3}) b + 4a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} + \right. \right. \\ & \quad \left. \left. (1+i\sqrt{3}) \left(-4a^3 + 18ab - 54c + 6\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{2/3} \right) \log \left[\left(2 \cdot 2^{1/3} (1-i\sqrt{3}) a^2 + 6i \cdot 2^{1/3} (1+\sqrt{3}) b + \right. \right. \\ & \quad \left. \left. 4a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} + (1+i\sqrt{3}) \left(-4a^3 + 18ab - 54c + 6\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{2/3} \right) \right] \right/ \\ & \quad \left. \left(12a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} \right) \right] \left/ \left(12a \left(-2a^3 + 9ab - 27c + 3\sqrt{3} \sqrt{-a^2 b^2 + 4b^3 + 4a^3 c - 18abc + 27c^2} \right)^{1/3} \log(3) \right) \right) \end{aligned}$$

就可以进一步得到 $t_1^2, t_2^2, t_1 t_2, t_1 \log(t_1), t_2 \log(t_2), t_1 \log(t_2), \log(t_1) \log(t_2), \dots$ 各种高次项和交叉项，在上述模型中选择了其中解释能力强的7个变量，实际上可以继续增添变量来提高预测能力，且复杂度不会变大。

八、总结

最开始做的时候，我钻进了牛角尖，试图用 a, b, c 具体地表达 H 并用泰勒展开进行级数化简，结果做了很久都没有成功。这时我突然想到，我应该先观察一下数据的特征，当我观察到 t_1, t_2, H 分布范围均在 $[0, 1]$ 之间，并通过画图发现的明显的线性相关性之后，一切都迎刃而解了，只需要利用简单的统计推断和线性回归的知识选择解释能力强的项优化模型即可。所以面对问题的时候一定要先观察，预处理和探索性分析非常重要！

总之，通过这次大作业，我对统计信号处理有了更加深刻的认识，对具体问题的分析过程得到了更多的锻炼，感谢老师和助教的耐心指导！

九、附录：源代码

matlab代码:

main.m

```
1  %%
2  clear all, close all, clc;
3
4  %% fit
5  %   using coefficient to fit H
6  %   t1 : b / a^2 (t1 > 0)
7  %   t2 : c / a^2 (t2 > 0)
8  %   t1t2 : cross term of t1 and t2
9  %   t1log : log(t1)
10 %   t2log : log(t2)
11 %   t1t1log : t1*log(t1)
12 %   t2t2log : t2*log(t2)
13 %   etc. Some terms may not be used.
14
15 train_num = 50000;
16 X = zeros(train_num, 8);
17 Y = zeros(train_num, 1);
18 for i = 1 : train_num
19     lambda = rand(3, 1) * 10000;
20     A = [lambda.^2, lambda, ones(3, 1)];
21     B = -lambda.^3;
22     t = (inv(A) * B)';
23     t = [t(2)/t(1)^2, t(3)/t(1)^3];
24     X(i, :) = [1, t, log(t), t.*log(t), prod(t)];
25     p = lambda / sum(lambda);
26     Y(i) = -p' * log(p) / log(3);
27 end
28
29 beta = inv(X'*X) * X' * Y;
30
31 %% test
32 n = 1000;
33 H = zeros(1, n);
34 Hhat = zeros(1, n);
35 %AX = B
36 %A : matrix of lambda_i lambda_i^2 1
37 %X : vector of a b c
38 %B : vector of lambda^3
39 for i = 1 : n
40     lambda = rand(3, 1) * 10000;
41     A = [lambda.^2, lambda, ones(3, 1)];
42     B = -lambda.^3;
43     t = (inv(A) * B)';
44     t = [t(2)/t(1)^2, t(3)/t(1)^3];
45     X = [1, t, log(t), t.*log(t), prod(t)];
46     p = lambda / sum(lambda);
47     H(i) = -p' * log(p) / log(3);
48     Hhat(i) = X * beta;
49 end
50
51 %% plot
52 %the estimate of H vs real H
53 error = mean(abs(Hhat - H));
54 plot(H, Hhat, 'o');
55 axis([0 1 0 1]);
56 xlabel('H');
57 ylabel('estimated H');
58 title('the estimate of H vs real H')
59 hold on;
60 x = 0 : 0.01 : 1;
61 y = 0 : 0.01 : 1;
62 plot(x, y);
63
64 %residual plot
65 figure();
66 index = 1:n;
67 plot(index, Hhat-H, 'o');
68 axis([0 n -0.01 0.01]);
69 xlabel('index');
70 ylabel('residuals');
71 title('residual plot')
72 hold on;
73 plot(index, index*0, 'r');
74
75 %% input
76 %input : lambda1,2,3
77 %output1 : a b c
```



```

78 %output2 : real H
79 %output3 : estimated H
80 %output4 : error
81 disp("You should have input coefficient a, b, c");
82 disp("but in this problem, a,b,c have strict conditions,");
83 disp("your input may lead to complex solutions, which is not allowed");
84 disp("instead, please input lambda1,2,3. I will calculate a,b,c according to your lambdas");
85 disp("and I will give the estimated H calculated from a,b,c, and the real H, and the error between them")
86 lambda1 = input("please input lambda1 which is positive:");
87 lambda2 = input("please input lambda2 which is positive:");
88 lambda3 = input("please input lambda3 which is positive:");
89 lambda = [lambda1, lambda2, lambda3]';
90 A = [lambda.^2, lambda, ones(3, 1)];
91 B = -lambda.^3;
92 t = (inv(A) * B)';
93 disp("a = "), disp(t(1)), disp("b = "), disp(t(2)), disp("c = "), disp(t(3));
94 t = [t(2)/t(1)^2, t(3)/t(1)^3];
95 X = [1, t, log(t), t.*log(t), prod(t)];
96 p = lambda / sum(lambda);
97 H = -p' * log(p) / log(3);
98 Hhat = X * beta;
99 disp("the real H equals "), disp(H);
100 disp("the estimated H equals "), disp(Hhat);
101 disp("the error is "), disp(Hhat-H);

```

R语言代码：

try.R

```

1 n = 10000
2 range = 10000
3 lambda1 = runif(n, 0, range)
4 lambda2 = runif(n, 0, range)
5 lambda3 = runif(n, 0, range)
6 A = array(0, dim = c(n, 3, 3))
7 X = array(0, dim = c(n, 3, 1))
8 B = array(0, dim = c(n, 3, 1))
9 Lambda = array(0, dim = c(n, 3, 1))
10 P = array(0, dim = c(n, 3))
11 H = array(0, dim = n)
12 for (i in 1:n) {
13     A[i, , ] = c(lambda1[i]^2, lambda2[i]^2, lambda3[i]^2,
14                   lambda1[i], lambda2[i], lambda3[i],
15                   1, 1, 1)
16     B[i, , ] = c(-lambda1[i]^3, -lambda2[i]^3, -lambda3[i]^3)
17     X[i, , ] = solve(A[i, , ])%%B[i, , ]
18     Lambda[i, , ] = c(lambda1[i], lambda2[i], lambda3[i])
19     P[i, ] = Lambda[i, , ] / (lambda1[i] + lambda2[i] + lambda3[i])
20     H[i] = -(P[i, 1]*logb(P[i, 1], 3) +
21             P[i, 2]*logb(P[i, 2], 3) +
22             P[i, 3]*logb(P[i, 3], 3))
23 }
24
25 a = X[, 1, ]
26 b = X[, 2, ]
27 c = X[, 3, ]
28 t1 = b/a^2
29 t2 = c/a^3
30 t1t2 = t1*t2
31 t1log = log(t1)
32 t2log = log(t2)
33 t1t1log = t1*t1log
34 t2t2log = t2*t2log
35 t1t2log = t1*t2log
36 t2t1log = t2*t1log
37 t1logt1log = t1log*t1log
38 t2logt2log = t2log*t2log
39 t1logt2log = t1log*t2log
40 fit = lm(H~t1+t2+t1log+t2log+t1t1log+t2t2log+t1t2)
41
42 error = mean(abs(fit$residuals))
43
44 plot(fit)

```

十、附录：文件清单

统计信号处理第二次大作业报告.pdf
main.m

大作业报告
matlab入口程序

