

732A90: Lab 4

Computational Statistics, Group 6

Chih-Yuan Lin, Sarah Walid Alsaadi, Carles Sans Fuentes, Joshua Burrata

March 3, 2017

Question 1

In this exercise, we are supposed to compute the Metropolis-Hastings algorithm from a target distribution which is

$$f(x) \propto x^5 * \exp^{-x}$$

. In order to evaluate whether we can get a good sample points from our algorithm, we are going to generate samples from the distribution by using as proposal distribution the $LN(X_t, 1)$, taking one starting point at our choice (in our case 1).

The explanation on how the data is sampled and simulated (from the algorithm) is as follows: (0) for as many points sampled we want, we are going to (1) sample a point from the proposed distribution (in our case from a lognormal with mean equal to X_t and sd equal to 1) which will be our y, and our point x_0 chosen at will. (2) Then, we generate a random uniform number equal to $U(0,1)$. (3) if

$$(\lognorm(y)/\lognorm(X_t)) * (p(X_t|y)/p(y|X_t)) > U(0, 1)$$

Then my new X_t value sampled is the previous X_t value else my new X_t is y. A plot representing 2000 random points sampled from a lognormal distribution using the Metropolis-Hastings can be found below in figure 1.

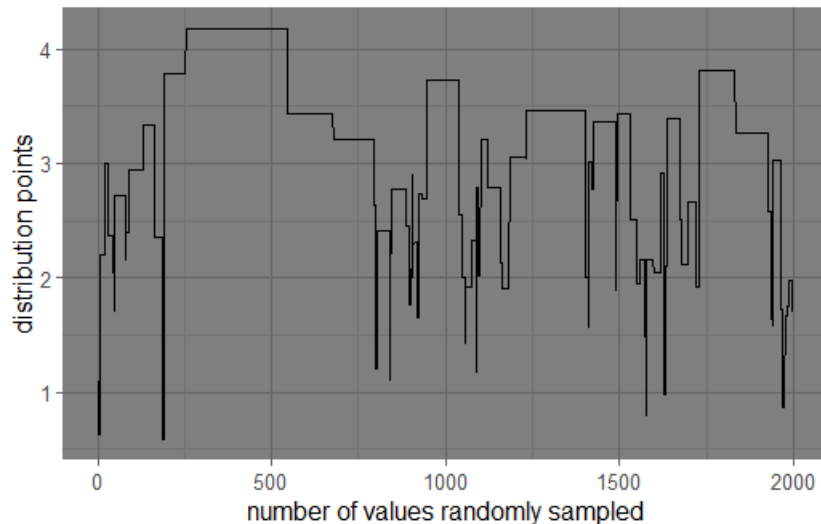


Figure 1: A plot representing 2000 random points sampled from a lognormal distribution using the Metropolis-Hastings

It can be seen that there is no convergence given the fact that the values get stuck so that there is no different sampled values chosen and this shows that our lognormal model is not good to

simulate our target distribution. Also, no burning period is clear since there is no convergence on the graphic.

Now, we are asked to do the same but sampling from a new proposed distribution being a chi-squared being $Chi(\lfloor X_t + 1 \rfloor)$. A plot representing 2000 random points sampled from a chi-squared distribution using the Metropolis-Hastings can be found below in figure 2.

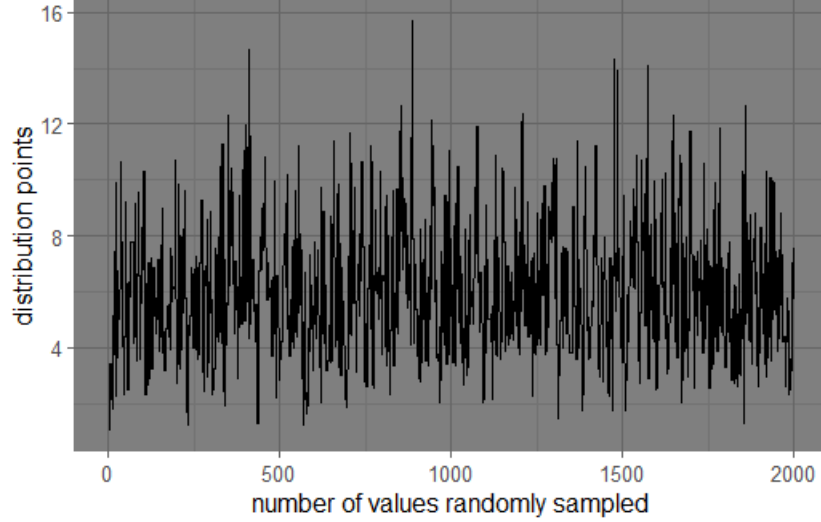


Figure 2: A plot representing 2000 random points sampled from a chi-squared distribution using the Metropolis-Hastings

It can be seen that there is convergence in this case given that our graphic looks totally as with a random walk. Nevertheless, a burning period is not clear to be found given the data. Using the Gelman-Rubin method, now we are trying to analyze the convergence of these chi-squared generated sequence. 10 MCMC sequences (from 1 to 10 by 1) have been generated using the "as.mcmc" function from the coda package, getting the following result:

```
1 attr(,"class")
2 [1] "mcmc.list"
3 > print(gelman.diag(mylist))
4 Potential scale reduction factors:
5
6      Point est. Upper C.I.
7 [1,]          1      1.01
```

The Gelman-Rubin method states that approximated convergence is diagnosed when the upper limit is close to 1. Therefore, we confirm convergence on our Metropolis-Hastings algorithm for our second Multi target distribution from which we simulate points. Nevertheless, the mean of our data is evaluated, being the means for our sampled distributions the following ones:

```
1 Elognormal
2 [1] 3.179977
3 > Echi
4 [1] 6.007571
```

The pdf of a gamma distribution with parameters (α, β) is:

$$f(x; \alpha, \beta) = \frac{\beta^\alpha x^{\alpha-1} e^{-x\beta}}{\Gamma(\alpha)}$$

The expectation of such a distribution is simply $E(X) = \alpha/\beta$.

Judging from the formulation of its pdf, the distribution we simulated here is clearly a gamma distribution with parameters $\alpha = 6, \beta = 1$. Therefore the expected value of this distribution and the actual value of the integral, is 6. So the estimated value of the integral from the step 2 sample is very close to this actual value, while the Step one sample estimate is much smaller.

Question 2

Exercise 2.1

In this exercise, we are asked to import some data to R and plot the dependence of Y on X. Here below in figure 3 it can be found the plot of the data:

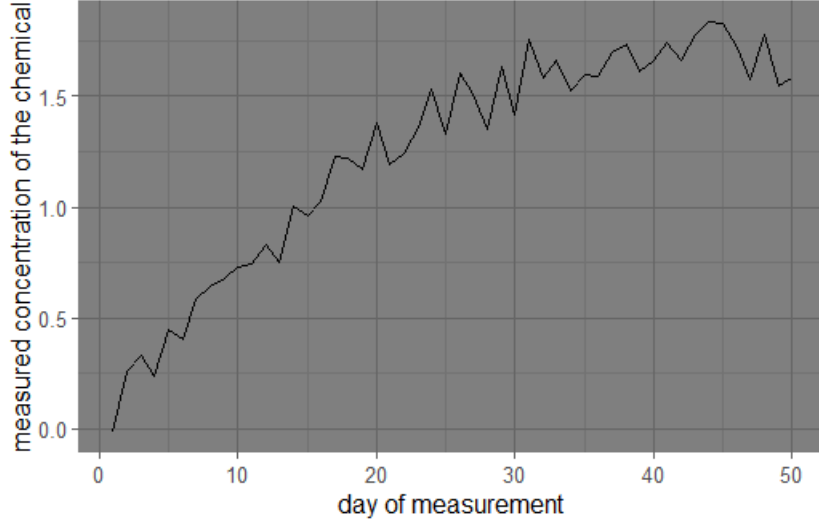


Figure 3: Plot the dependence of Y on X

The data seems to follow a logarithmic distribution.

Exercise 2.2

The likelihood is given by:

$$p(Y|\mu) = \prod_{i=1}^n p(y_i|\mu) \propto \exp\left(-\frac{1}{2 * 0.2^2} \sum_{i=1}^n (y_i - \mu_i)^2\right)$$

and the prior is given by:

$$p(\mu) = \prod_{i=1}^n p(\mu_i) \propto \exp\left(-\frac{1}{2 * 0.2^2} \sum_{i=2}^n (\mu_i - \mu_{i-1})^2\right)$$

Thus, the posterior is given by:

$$p(\mu|Y) \propto \exp\left(-\frac{1}{2 * 0.2^2} \left(\sum_{i=1}^n (y_i - \mu_i)^2 + \sum_{i=2}^n (\mu_i - \mu_{i-1})^2\right)\right)$$

.

Exercise 2.3

From the previous exercise, we get the full conditional distribution as:

part 1:

$$p(\mu_1|\mu_{-1}, Y) \propto \exp\left(-\frac{1}{2 * 0.2^2} ((\mu_1 - y_1)^2 + (\mu_1 - \mu_2)^2)\right) \propto \exp\left(-\frac{1}{\frac{2 * 0.2^2}{2}} \left(\mu_1 - \frac{(y_1 + \mu_2)}{2}\right)^2\right)$$

part 2:

$$p(\mu_i | \mu_{-i}, Y) \propto \exp\left(-\frac{1}{2 * 0.2^2} ((\mu_i - y_i)^2 + (\mu_i - \mu_{i-1})^2 + (\mu_i - \mu_{i+1})^2)\right) \propto \exp\left(-\frac{1}{\frac{2 * 0.2^2}{3}} \left(\mu_i - \frac{(y_i + \mu_{i-1} + \mu_{i+1})}{3}\right)^2\right)$$

for $i = 2, 3, \dots, n - 1$.

part 3:

$$p(\mu_n | \mu_{-n}, Y) \propto \exp\left(-\frac{1}{2 * 0.2^2} ((\mu_n - y_n)^2 + (\mu_n - \mu_{n-1})^2)\right) \propto \exp\left(-\frac{1}{\frac{2 * 0.2^2}{2}} \left(\mu_n - \frac{(y_n + \mu_{n-1})}{2}\right)^2\right)$$

Exercise 2.4

We then used 3 previously derived expressions to implement a Gibbs sampler with starting point $\vec{\mu}^0$. We ran the Gibbs sampler to obtain 1000 $\vec{\mu}$ values. We took the mean over these 1000 for each component of $\vec{\mu}$ and plotted them against X, alongside Y vs X, as shown in Figure 4

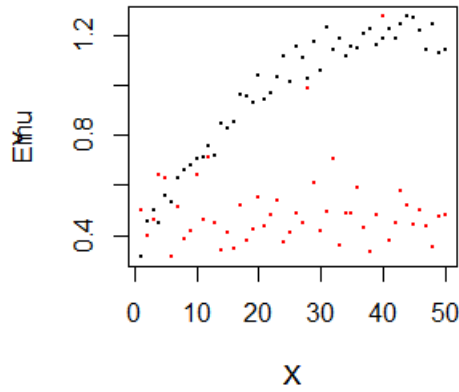


Figure 4: Plot showing Y (in black) and the expected Mu values (in red) vs X

Exercise 2.5

Figure 5 is a trace plot of μ_n over its 1000 iterations.

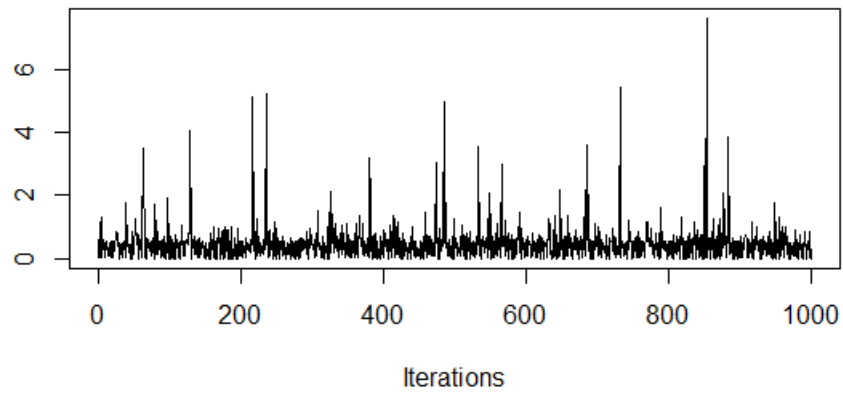


Figure 5: Trace Plot of μ for $i=n$

Comment: The results in 2.4 and 2.5 seem a little strange, as μ is not close to Y at all, so we think that our Gibbs sampler is incorrect.

Contributions

All results and comments presented have been developed and discussed together by the members of the group.

Appendix

Assignment 1

```
1  ###LAB 4
2  ## Activity 1
3
4  set.seed(12345)
5  ###1.1
6  posterior_PDFlognorm<- function(value,param){
7    return(dlnorm(x=value, meanlog = param, sdlog = 1))
8  }
9
10 ##target distribution
11 targetDist<- function(x){
12   return(x**5*exp(-x))
13 }
14
15 MetropolisHastings<- function(times, init, posterior,ydist="rlnorm"){
16   result<- integer(0)
17   result[1]<-init
18   counter<-1
19   init<- init
20
21   while(counter<times){
22
23     x<-init
24     if(ydist == "chi"){
25       y0<- rchisq(1,floor(x+1))
26     }
27     else{
28       y0<- rlnorm(1, x,1)
29     }
30     #Sample a candidate value X* from a proposal distribution g(.|x)
31
32     myrunif<- runif(1,0,1)
33
34     alpha<-(posterior(x, y0)/posterior(y0,x))*(targetDist(y0)/targetDist(x))
35     if(myrunif<alpha){
36       init<- y0
37     }
38     else{
39       init<-init
40     }
41     counter<-counter+1
42     result[counter]<- init
43   }
44   return(result)
45 }
46
47 MCMH<-MetropolisHastings(2000,1, posterior_PDFlognorm)
48
49 library(ggplot2)
50 myplot<- function(data){
51   mydata<- as.data.frame(data)
52   mplot<- ggplot2::ggplot(data = mydata)+geom_line(aes( x= 1:nrow(mydata), y = data))+
53     ylab("distribution points")+xlab("number of values randomly sampled")+
54     theme_dark()
55   return(mplot)
56 }
57
58 myplot(MCMH)
59
60
61 ###1.2
62 set.seed(12345)
63 chisquarePDF<- function(value,param){
64   return(dchisq(value,floor(param+1)))
65 }
66
67 MCMH_chi<-MetropolisHastings(2000, 1, chisquarePDF, "chi")
68 myplot(MCMH_chi)
69
70 ###1.3 EXplanation
71
72 ###1.4
73 #install.packages("coda")
74 library(coda)
75 mylist<-mcmc.list()
76
77 times<-2000
78 n<- 10
79
80 for( i in 1:n){
81   mylist[[i]]<-as.mcmc(MetropolisHastings(times, i, chisquarePDF, "chi"))
82 }
```

```

82 }
83 }
84 mylist
85 print(gelman.diag(mylist))
86
87
88 ##1.5
89 Elognormal<- mean(MCMH)
90 Echi<- mean(MCMH_chi)
91
92 ##1.6
93 alpha <- 6
94 beta<- 1
95 Ebeta <- alpha/Beta

```

Assignment 2

```

1 ##Activity 2
2 library(ggplot2)
3 load("C:/Users/M/Desktop/Statistics and Data Mining Master/Semester 2/Computational Statistics/
  Lecture 4/chemical.RData")
4
5 ##2.1
6 mydata<- data.frame(X=X, Y = Y)
7
8 ggplot2::ggplot(mydata)+ geom_line(aes(x = X, y = Y))+
9   xlab("day of measurement")+ ylab("measured concentration of the chemical")+
10   theme_dark()
11
12 ##2.4
13 gibbs <- function(tmax, mu0, var=0.2) {
14   d <- length(mu0)
15   mu <- matrix(0, nrow=tmax, ncol = d)
16   mu[1, ] <- mu0
17   Y <- matrix(0, nrow=tmax, ncol = d)
18   Y[1, ] <- sapply(X=mu0, FUN=function(mu){
19     y <- rnorm(1, mean=mu, sd=0.2)
20     return(y)
21   })
22   t = 1
23   while (t < tmax) {
24     #update mu_t
25     #i=1
26     mu[t+1, 1] <- exp(-1/0.04*(mu[t, 1]-(Y[t, 1]+mu[t, 2])/2)^2)/(dnorm(Y[t, 1], mean=mu[t, 1],
27       sd=0.2))
28     #i =2,3,...,d-1
29     dmin1 <- d-1
30     for (i in 2:dmin1) {
31       mu[t+1, i] <- exp(-1/(0.08/3)*(mu[t, i]-(Y[t, i]+mu[t+1, i-1]+mu[t, i+1])/3)^2)/(dnorm(Y[
32         t, i], mean=mu[t, i], sd=0.2))
33     }
34     #i=d
35     mu[t, d] <- exp(-1/0.04*(mu[t, d]-(Y[t, d]+mu[t+1, d-1])/2)^2)/(dnorm(Y[t, d], mean=mu[t, d
36       ], sd=0.2))
37
38     #update Y_t
39     Y[t+1, ] <- sapply(X=mu[t+1, ], FUN=function(mu){
40       y <- rnorm(1, mean=mu, sd=0.2)
41       return(y)
42     })
43
44     t <- t+1
45   }
46   return(mu)
47 }
48
49 mu0 <- rep(0, 50)
50 tmax <- 1000
51 mu <- gibbs(tmax=tmax, mu0=mu0)
52 Emu <- colMeans(mu)
53 plot(X, Emu,pch=19, cex=0.3, col="red")
54 par(new=TRUE)
55 plot(X, Y,pch=19, cex=0.3, col="black", axes = FALSE)
56
57 ##2.5
58 mun <- as.mcmc(mu[, 50])
59 library(coda)
60 traceplot(mun, type="p", pch=19, cex=0.3)

```