

Bayesian Lab3

Xuan Wang & Lepeng Zhang

2024-05-13

Question 1 Gibbs sampling for the logistic regression

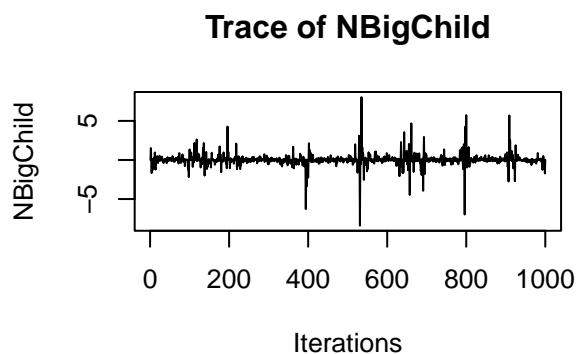
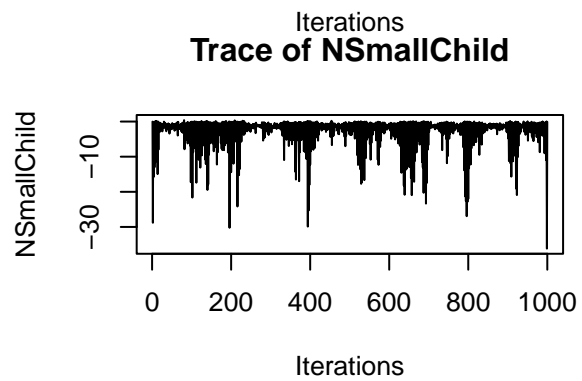
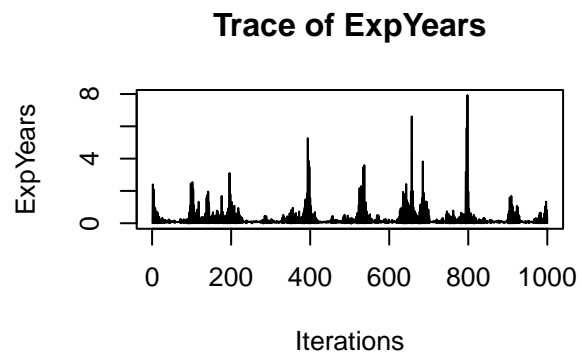
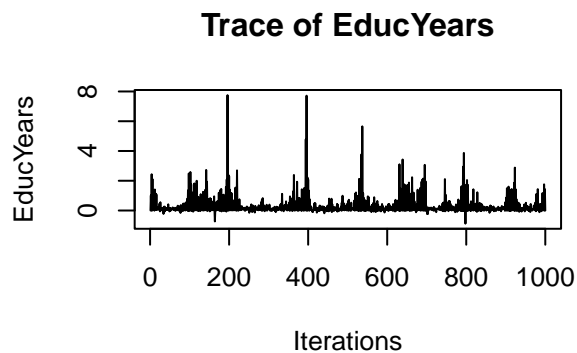
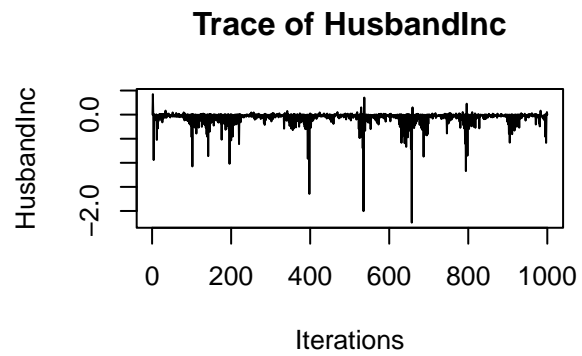
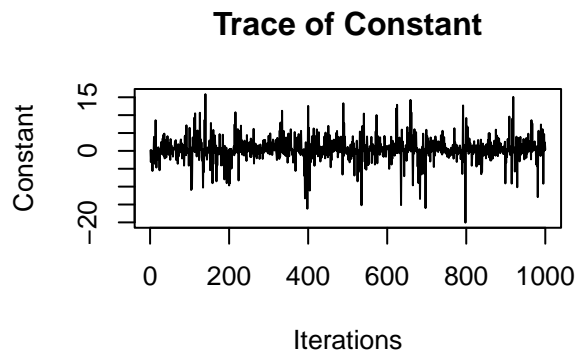
1a)

```
library("mvtnorm")
library("BayesLogit")
set.seed(12345)
raw_data <- read.table("WomenAtWork.dat", header= TRUE)
y <- raw_data$Work
X <- as.matrix(raw_data[,2:ncol(raw_data)])
Xnames <- colnames(X)

n <- dim(X)[1]
Npar <- dim(X)[2]
tau <- 3
b <- as.matrix(rep(0,Npar))
B <- (tau^2)*diag(Npar)
k <- y-0.5
init_beta <- rmvnorm(1, mean = b, sigma = B)
n_iter <- 1000
w <- rep(0,n)
post_beta <-matrix(0, nrow= n_iter, ncol= Npar)

for (i in 1:n_iter){
  for (j in 1:n){
    a <- X[j,]%*%t(init_beta)
    w[j] <- rpg(1, abs(a))
  }
  O <- diag(w)
  V <- solve(t(X)%*%O%*%X + solve(B))
  m <- V%*%(t(X)%*%k + solve(B)%*%b)
  beta <- rmvnorm(1, m, V)
  init_beta <- beta
  post_beta[i,] <- beta
}

par(mfrow= c(2, 2))
for (i in 1:Npar){
  plot(post_beta[,i], main= paste("Trace of",Xnames[i]), type= "l", xlab='Iterations', ylab=Xnames[i])
}
```



```
IFs <- apply(post_beta, 2, function(x){
  acf_x <- acf(x, plot = FALSE)$acf[,1]
  return(1 + 2 * sum(acf_x[-1]))
})
names(IFs) <- Xnames
```

IFs

```
##      Constant  HusbandInc  EducYears  ExpYears      Age NSmallChild
##      1.051941    3.434646    3.741895    4.253009    2.737143    4.296229
##      NBigChild
##      0.466449
```

In the terms of IFs and trace plots, the sampling has not converged well except for Constant and NBigChild as their IFs are less than 1.1.

1b)

```
new_x <- c(1, 22, 12, 7, 38, 1, 0)
logistic <- function(x){
  1/(1+exp(-x))
}
prob <- logistic(new_x %*% t(post_beta))
quantile(prob, c(0.05, 0.95))
```

```
##          5%          95%
## 3.694238e-05 5.061599e-01
```

```
#hist(prob, freq=FALSE, breaks = 500)
```

2a)

```
set.seed(12345)
raw_data <- read.table("eBayNumberOfBidderData_2024.dat", header = TRUE)
y <- raw_data[,1]
X <- as.matrix(raw_data[,-1])
Xnames <- colnames(X)
m1 <- glm(y~X[, -1], raw_data, family='poisson')
summary(m1)
```

```
##
## Call:
## glm(formula = y ~ X[, -1], family = "poisson", data = raw_data)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.07981    0.03393  31.828 < 2e-16 ***
## X[, -1]PowerSeller -0.03566    0.04167  -0.856 0.392109
## X[, -1]VerifyID   -0.45564    0.12748  -3.574 0.000351 ***
## X[, -1]Sealed      0.45515    0.06226   7.311 2.65e-13 ***
## X[, -1]Minblem    -0.06837    0.07198  -0.950 0.342228
## X[, -1]MajBlem    -0.22554    0.09525  -2.368 0.017894 *
## X[, -1]LargNeg     0.05382    0.06406   0.840 0.400787
## X[, -1]LogBook    -0.08499    0.03234  -2.628 0.008599 **
## X[, -1]MinBidShare -1.82490    0.07843 -23.269 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
```

```
##      Null deviance: 1699.6  on 799  degrees of freedom
## Residual deviance:  691.8  on 791  degrees of freedom
## AIC: 2879.1
##
## Number of Fisher Scoring iterations: 5
```

VerifyID, Sealed, MinBidShare are important.

2b)

```
n <- dim(X)[1]
p <- dim(X)[2]
mu <- rep(0, p)
Sigma <- 100*solve(t(X)%*%X)

log_posterior <- function(beta,X,y){
  log_likelihood = sum(y*(X %*% beta)) - sum(exp(X %*% beta))
  log_prior = dmvnorm(beta, mu, Sigma, log = TRUE)
  return(log_likelihood + log_prior)
}

optim_result <- optim(coef(m1), log_posterior, gr=NULL, X, y, method=c('BFGS'), control=list(fnscale=-1))

beta_hat <- optim_result$par
names(beta_hat) <- Xnames
print('The posterior modes:')

## [1] "The posterior modes:"
print(beta_hat)

##      Const PowerSeller  VerifyID      Sealed      Minblem      MajBlem
## 1.07724892 -0.03571573 -0.45299504  0.45482368 -0.06861145 -0.22597142
##      LargNeg      LogBook MinBidShare
## 0.05398896 -0.08452628 -1.82252976

J <- -optim_result$hessian
inv_J <- solve(J)
rownames(inv_J) <- Xnames
colnames(inv_J) <- Xnames
print('The posterior covariance:')

## [1] "The posterior covariance:"
print(inv_J)

##      Const      PowerSeller      VerifyID      Sealed
## Const      1.148876e-03 -8.900221e-04 -0.0003857695 -3.819207e-04
## PowerSeller -8.900221e-04  1.736885e-03 -0.0001082671 -3.044257e-04
## VerifyID    -3.857695e-04 -1.082671e-04  0.0161606433 -9.387069e-04
## Sealed      -3.819207e-04 -3.044257e-04 -0.0009387069  3.877766e-03
## Minblem     -5.338585e-04  7.431093e-05  0.0001663128  4.467620e-04
## MajBlem     -3.275944e-04 -2.789405e-04  0.0003412998  5.282208e-04
## LargNeg     -6.095215e-04  3.646050e-04  0.0003556882  3.758890e-04
## LogBook      4.158001e-05  1.732838e-04 -0.0003756359 -5.804635e-05
## MinBidShare 1.323796e-03 -6.728060e-04 -0.0008294539 -1.321082e-04
##      Minblem      MajBlem      LargNeg      LogBook
```

```
## Const      -5.338585e-04 -0.0003275944 -6.095215e-04  4.158001e-05
## PowerSeller 7.431093e-05 -0.0002789405  3.646050e-04  1.732838e-04
## VerifyID    1.663128e-04  0.0003412998  3.556882e-04 -3.756359e-04
## Sealed      4.467620e-04  0.0005282208  3.758890e-04 -5.804635e-05
## Minblem     5.181453e-03  0.0004407869  6.458399e-05 -1.390083e-06
## MajBlem     4.407869e-04  0.0090786284  5.029311e-04 -1.357592e-04
## LargNeg     6.458399e-05  0.0005029311  4.106042e-03 -3.194993e-04
## LogBook     -1.390083e-06 -0.0001357592 -3.194993e-04  1.045587e-03
## MinBidShare -1.987332e-04  0.0002878482 -5.369607e-05  1.247716e-03
##           MinBidShare
## Const      1.323796e-03
## PowerSeller -6.728060e-04
## VerifyID    -8.294539e-04
## Sealed      -1.321082e-04
## Minblem     -1.987332e-04
## MajBlem     2.878482e-04
## LargNeg     -5.369607e-05
## LogBook     1.247716e-03
## MinBidShare 6.125572e-03
```

The posterior mode is close to the coefficients in 1a).

2c)

```
logPostFunc_poi <- function(theta, X, y){
  log_likelihood = sum(y*(X %*% theta)) - sum(exp(X %*% theta))
  log_prior = dmvnorm(theta, mu, Sigma, log = TRUE)
  return(log_likelihood + log_prior)
}

RWMSampler <- function(Func, c_value, init_theta, X, y, n_samples=10000){
  return_list <- list()
  draws <- matrix(0, n_samples, ncol(X))
  i <- 1
  draws[i,] <- init_theta
  count_accept <- 0

  while(i<n_samples){
    i = i + 1
    theta_p <- as.numeric(rmvnorm(1, draws[i-1,], c_value*inv_J))
    acc_pro <- min(1, exp(Func(theta_p, X, y)-Func(draws[i-1,], X, y)))

    if(runif(1)<acc_pro){
      draws[i,] <- theta_p
      count_accept <- count_accept+1
    } else{
      draws[i,] <- draws[i-1,]
    }
  }
  accept_rate <- count_accept/(n_samples-1)
  return_list$draws <- draws
  return_list$accept_rate <- accept_rate
  return(return_list)
}
```

```

set.seed(12345)
theta_0 <- rep(0, ncol(X))
result <- RWMSampler(logPostFunc_poi, c_value = 0.5, init_theta = theta_0, X = X, y = y)
draws <- result$draws
accept_rate <- result$accept_rate

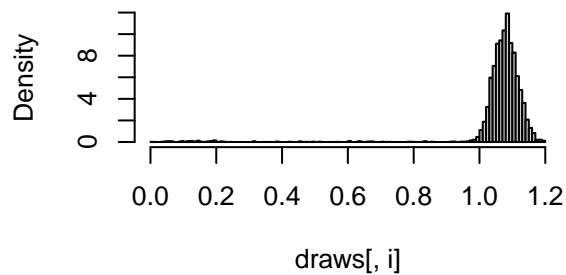
cat(paste0('The acceptance probability is ',round(accept_rate*100,1), '%.'))

## The acceptance probability is 32.2%.

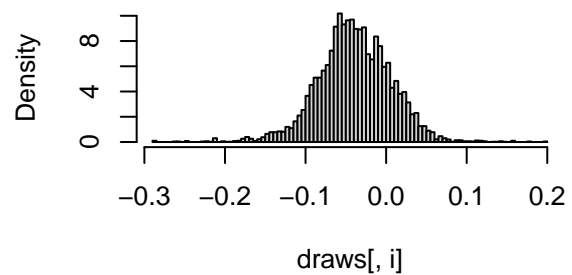
par(mfrow = c(2, 2))
for (i in 1:p) {
  hist(draws[,i], main = paste("Hist of", Xnames[i]), breaks = 100, freq = FALSE)
}

```

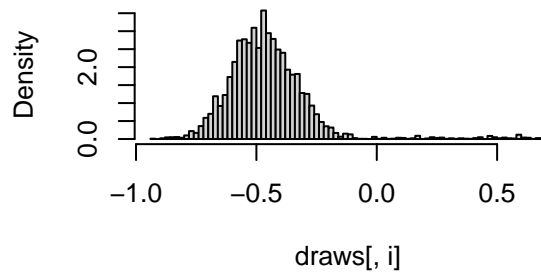
Hist of Const



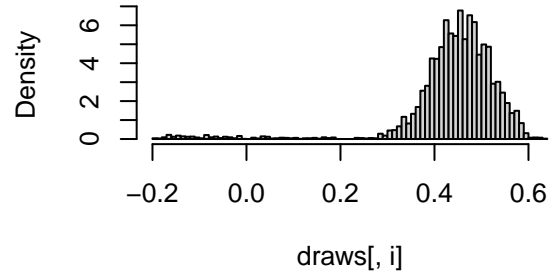
Hist of PowerSeller



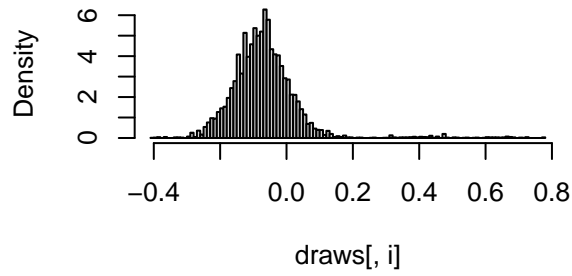
Hist of VerifyID



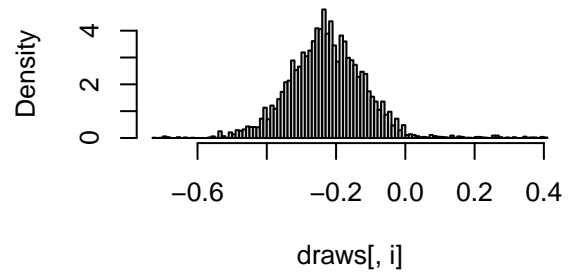
Hist of Sealed



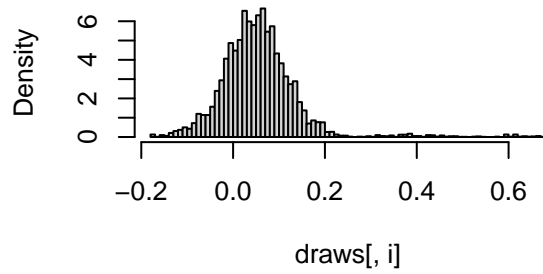
Hist of Minblem



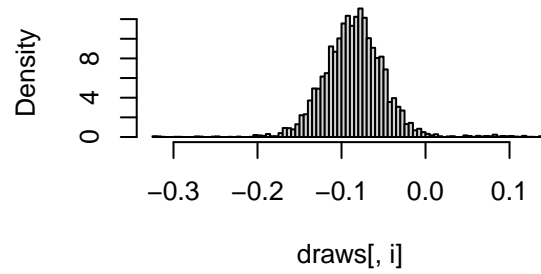
Hist of MajBlem



Hist of LargNeg

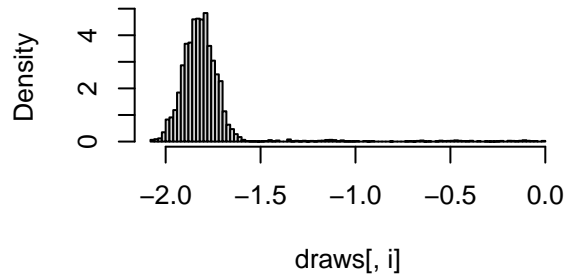


Hist of LogBook

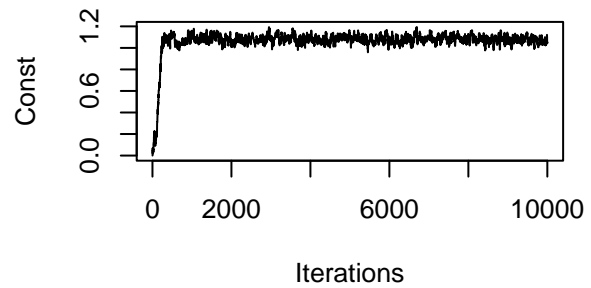


```
for (i in 1:p) {  
  plot(draws[, i], main = paste("Trace of", Xnames[i]), type = "l", xlab='Iterations', ylab= Xnames[i])  
}
```

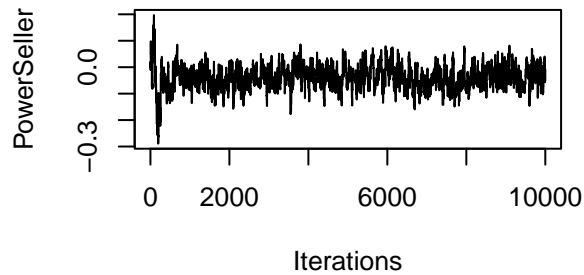
Hist of MinBidShare



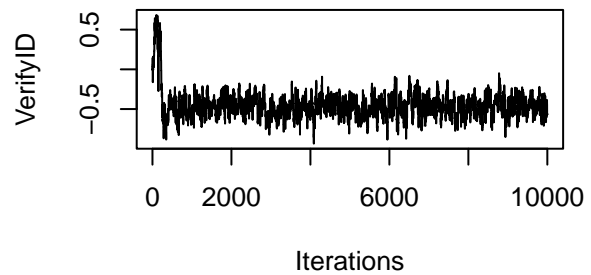
Trace of Const

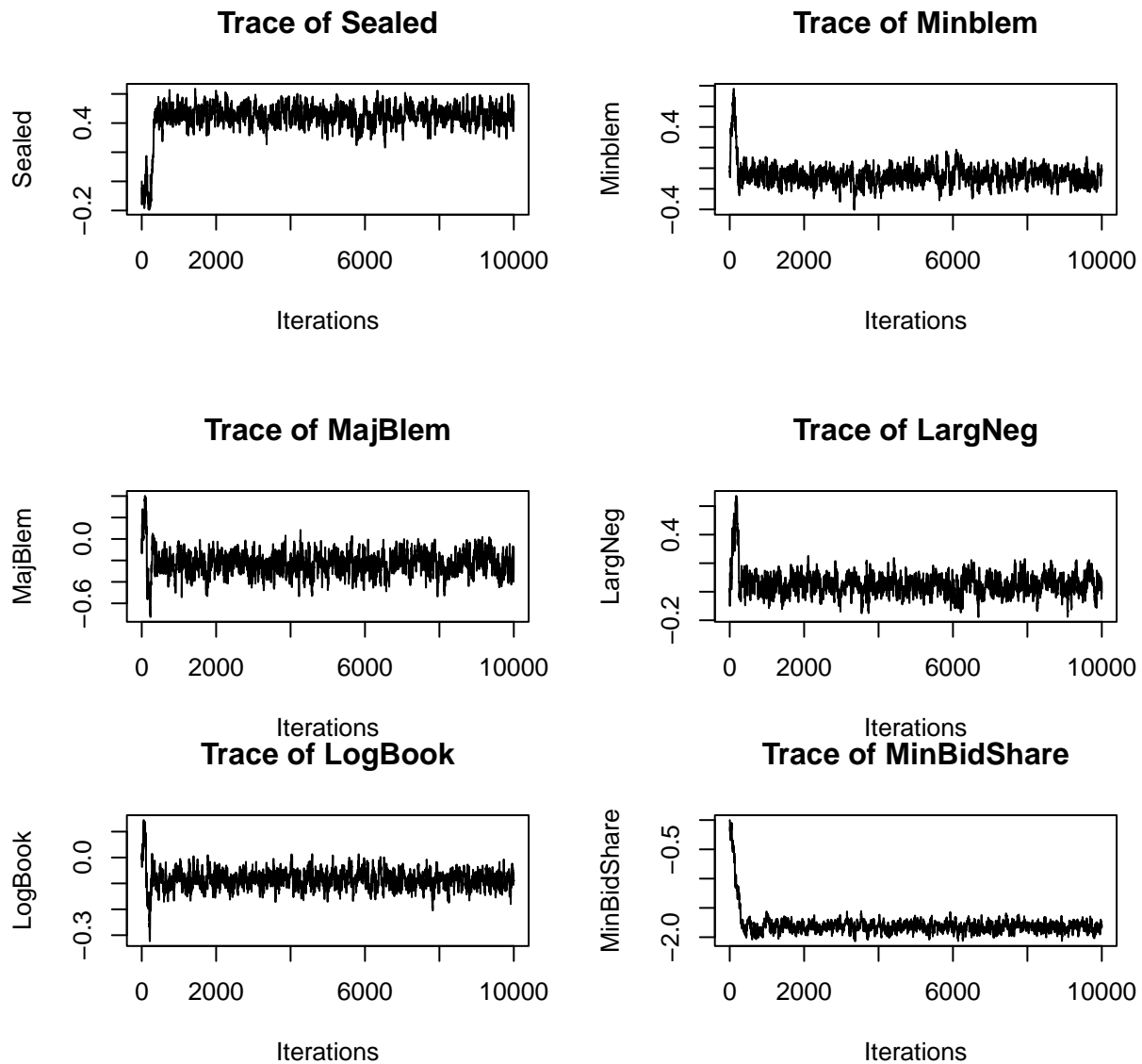


Trace of PowerSeller



Trace of VerifyID





From both histograms and trace plots, we can see that they all converged.

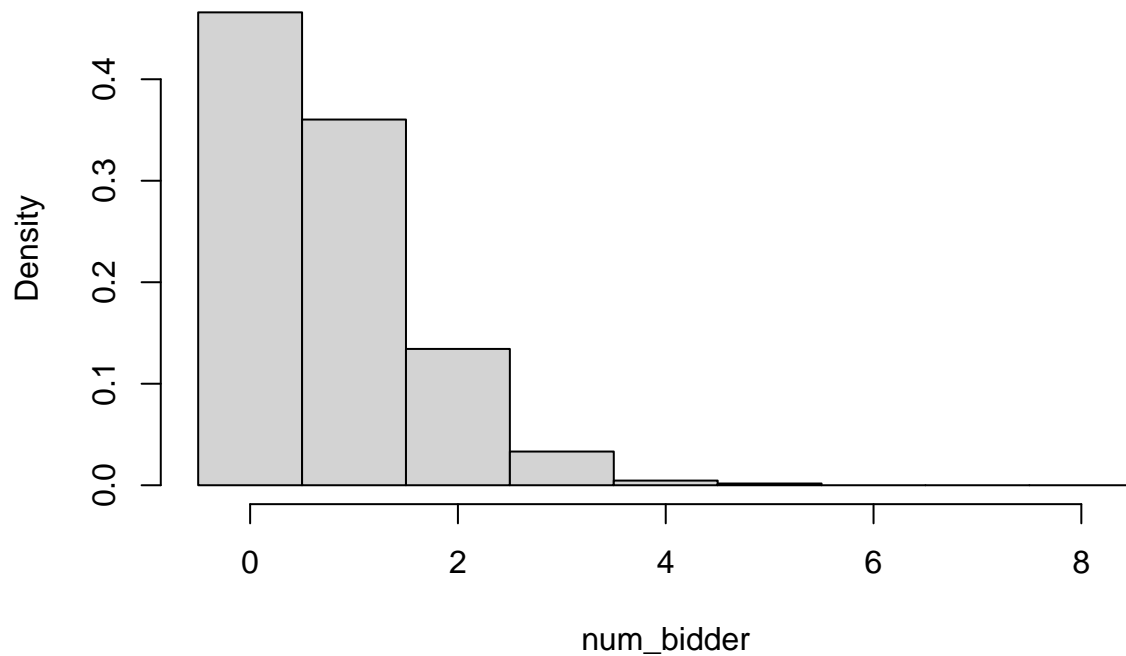
2d)

```
X_new <- c(1, 1, 0, 1, 0, 1, 0, 1.2, 0.8)

num_bidder <- c()
for (i in 1:nrow(draws)){
  num_bidder[i] <- rpois(1, lambda = exp(draws[i,] %*% X_new))
}

hist(num_bidder, breaks = seq(-0.5, 8.5, 1), freq = FALSE)
```


Histogram of num_bidder



```
cat(paste0('The prob that no bidders is ',round(mean(num_bidder==0),2),'.'))
```

```
## The prob that no bidders is 0.47.
```

Question 3 Time series models in Stan

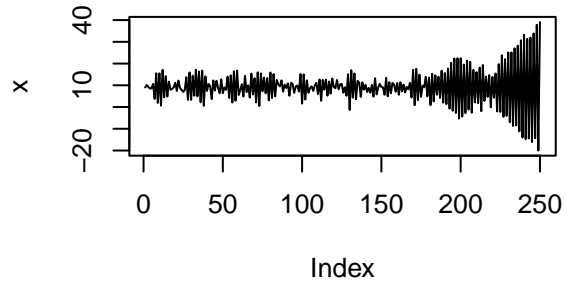
3a)

```
# Load required library
library(rstan)
set.seed(12345)
# Function to simulate AR(1) process
simulate_ar1 <- function(mu, phi, sigma_sq, T) {
  x <- numeric(T)
  x[1] <- mu
  for (t in 2:T) {
    e_t <- rnorm(1, mean = 0, sd = sqrt(sigma_sq))
    x[t] <- mu + phi * (x[t-1] - mu) + e_t
  }
  return(x)
}

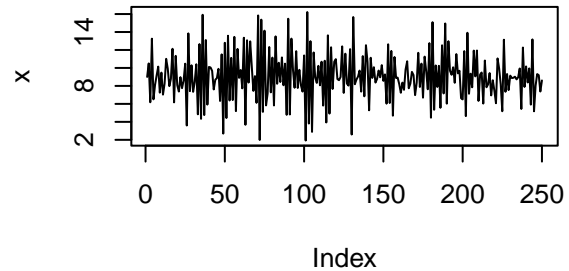
# Simulate AR(1) process
mu <- 9
sigma_sq <- 4
T <- 250
phi_values <- seq(-1, 1, by = 0.2)
par(mfrow = c(2, 2))
for (phi in phi_values) {
```

```
x <- simulate_ar1(mu, phi, sigma_sq, T)
plot(x, main = paste("AR(1) process with phi =", phi), type = "l")
}
```

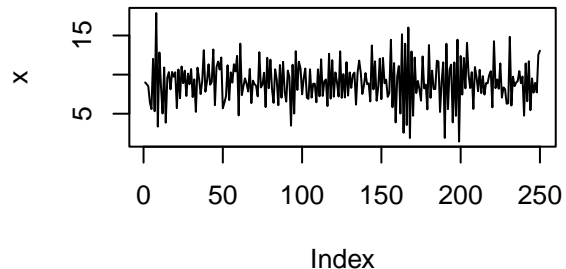
AR(1) process with $\phi = -1$



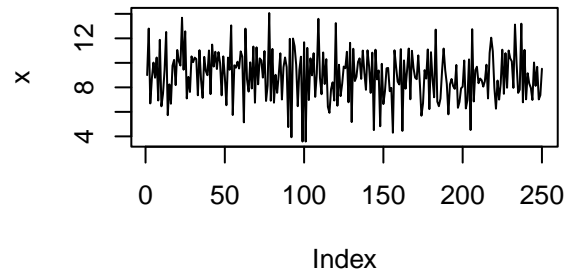
AR(1) process with $\phi = -0.8$



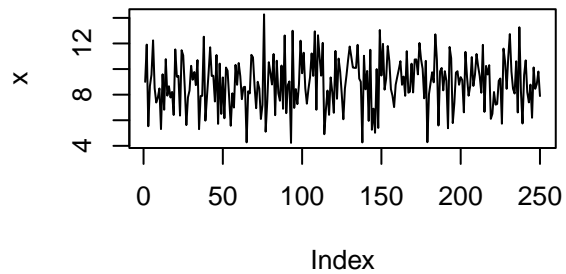
AR(1) process with $\phi = -0.6$



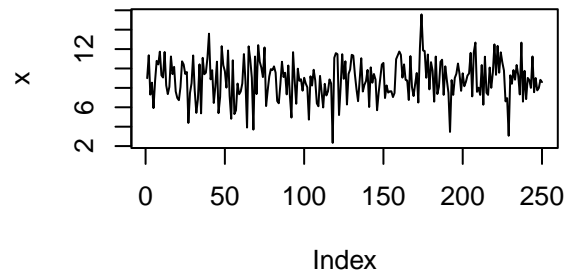
AR(1) process with $\phi = -0.4$



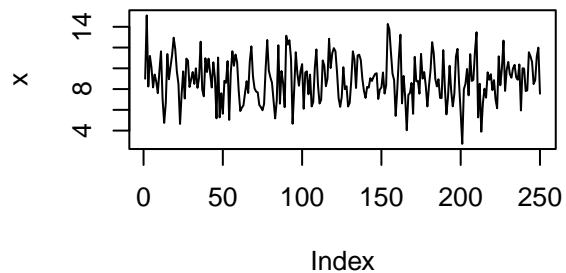
AR(1) process with $\phi = -0.2$



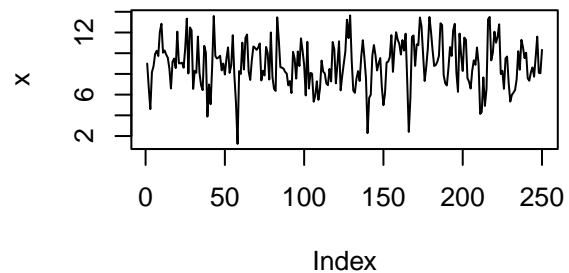
AR(1) process with $\phi = 0$



AR(1) process with $\phi = 0.2$

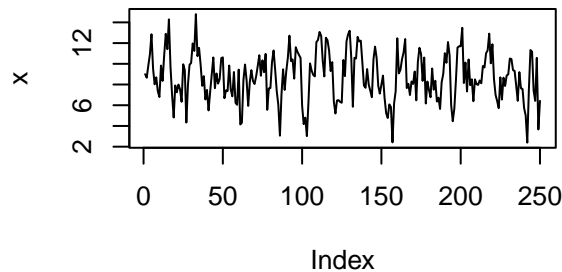


AR(1) process with $\phi = 0.4$

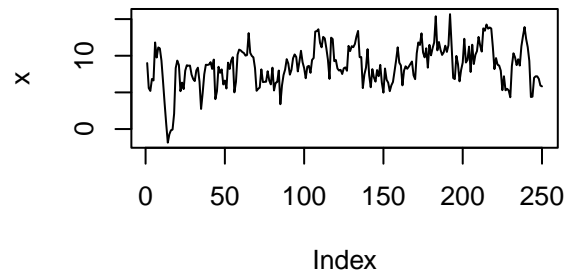


```
par(mfrow = c(1, 1)) # Reset the plot layout
```

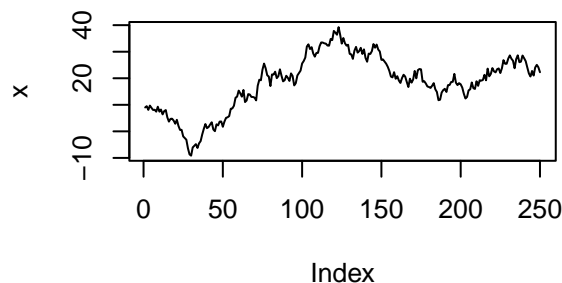
AR(1) process with $\phi = 0.6$



AR(1) process with $\phi = 0.8$



AR(1) process with $\phi = 1$



ϕ close to -1: The plot shows a strong negative autocorrelation. That is, if the process was high (or low) at one point, it is likely to be low (or high) at the next point. This aligns with the first plot that there is a wider up and down range as the process frequently changes between high and low values.

ϕ close to 0: The plot shows there are not too much autocorrelation. The values in the process are not influenced by their preceding values. As a result, it can be seen that the fluctuation becomes smooth and within a regular range, meaning the process appears more random and less predictable.

ϕ close to 1: The plot shows a strong positive autocorrelation. If the process was high (or low) at one point, it is likely to be high (or low) at the next point. This can lead to irregular up and down fluctuations as the process tends to maintain its current state for longer periods, as the last plot shown.

```
set.seed(12345)
# Stan code
stan_code <- "
data {
  int<lower=0> T;
  vector[T] y;
}
parameters {
  real mu;
  real<lower=-1, upper=1> phi;
  real<lower=0> sigma;
}
model {
  vector[T] nu;
  nu[1] <- mu;
  for (t in 2:T) {
```

```

    nu[t] <- mu + phi * (y[t-1] - mu);
  }
  y ~ normal(nu, sigma);
}
"

# Compile Stan model
stan_model <- stan_model(model_code = stan_code)

# Simulate two AR(1) processes
x1 <- simulate_ar1(mu, 0.3, sigma_sq, T)
y1 <- simulate_ar1(mu, 0.97, sigma_sq, T)

# Fit Stan model to the simulated data
fit_x1 <- sampling(stan_model, data = list(T = T, y = x1))
fit_y1 <- sampling(stan_model, data = list(T = T, y = y1))

# Print the results
print(summary(fit_x1)$summary)

##                mean      se_mean      sd      2.5%      25%      50%
## mu          9.1905497 0.002709162 0.15520180  8.895450  9.0862328  9.1918514
## phi         0.2424291 0.001016187 0.06148230  0.125924  0.2006003  0.2421251
## sigma       1.8490849 0.001317939 0.08381263  1.689257  1.7903087  1.8469382
## lp__       -278.1752514 0.026466853 1.21964442 -281.273451 -278.7220575 -277.8755240
##                75%      97.5%    n_eff    Rhat
## mu          9.2957065   9.4892556 3281.886 1.0011961
## phi         0.2838876   0.3694193 3660.603 0.9995990
## sigma       1.9031192   2.0221761 4044.160 0.9993457
## lp__       -277.2639605 -276.7786754 2123.547 1.0018915

print(summary(fit_y1)$summary)

##                mean      se_mean      sd      2.5%      25%
## mu          8.0753019 0.036446688 1.82916093  4.5986100  6.8420866
## phi         0.9695688 0.000362634 0.01573012  0.9355659  0.9595571
## sigma       1.9141773 0.001706622 0.08598444  1.7495130  1.8578694
## lp__       -289.4189068 0.043535591 1.39998956 -293.1693479 -290.0527891
##                50%      75%      97.5%    n_eff    Rhat
## mu          7.9923555   9.2803634 11.9913829 2518.765 1.0011927
## phi         0.9706627   0.9811544  0.9964299 1881.598 0.9999003
## sigma       1.9097000   1.9676696  2.0933737 2538.430 0.9997161
## lp__       -289.0323030 -288.3924958 -287.8628621 1034.096 1.0020259

```

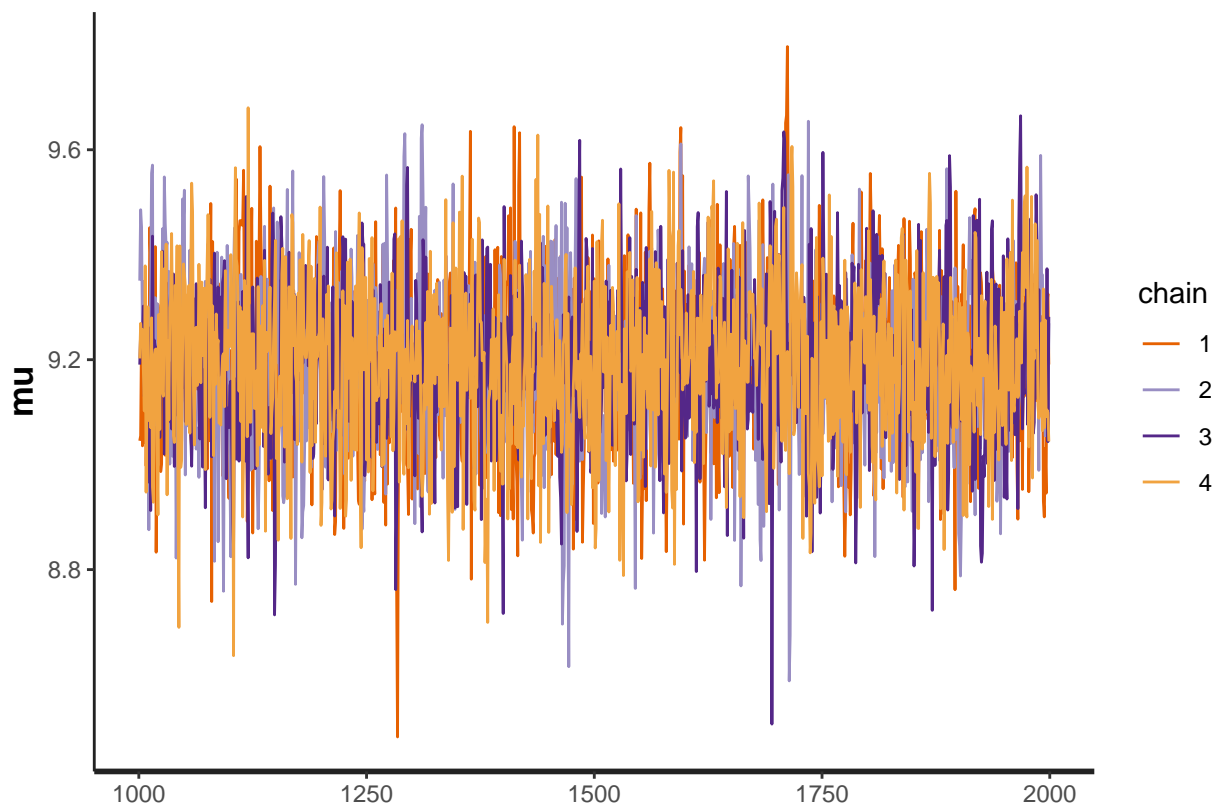
From above results, for the first dataset from $\phi=0.3$,

parameters	mean	95% credible intervals	the number of effective samples
mu	9.191	8.895-9.489	3281
phi	0.242	0.201-0.369	3660
sigma	1.849	1.689-2.022	4044

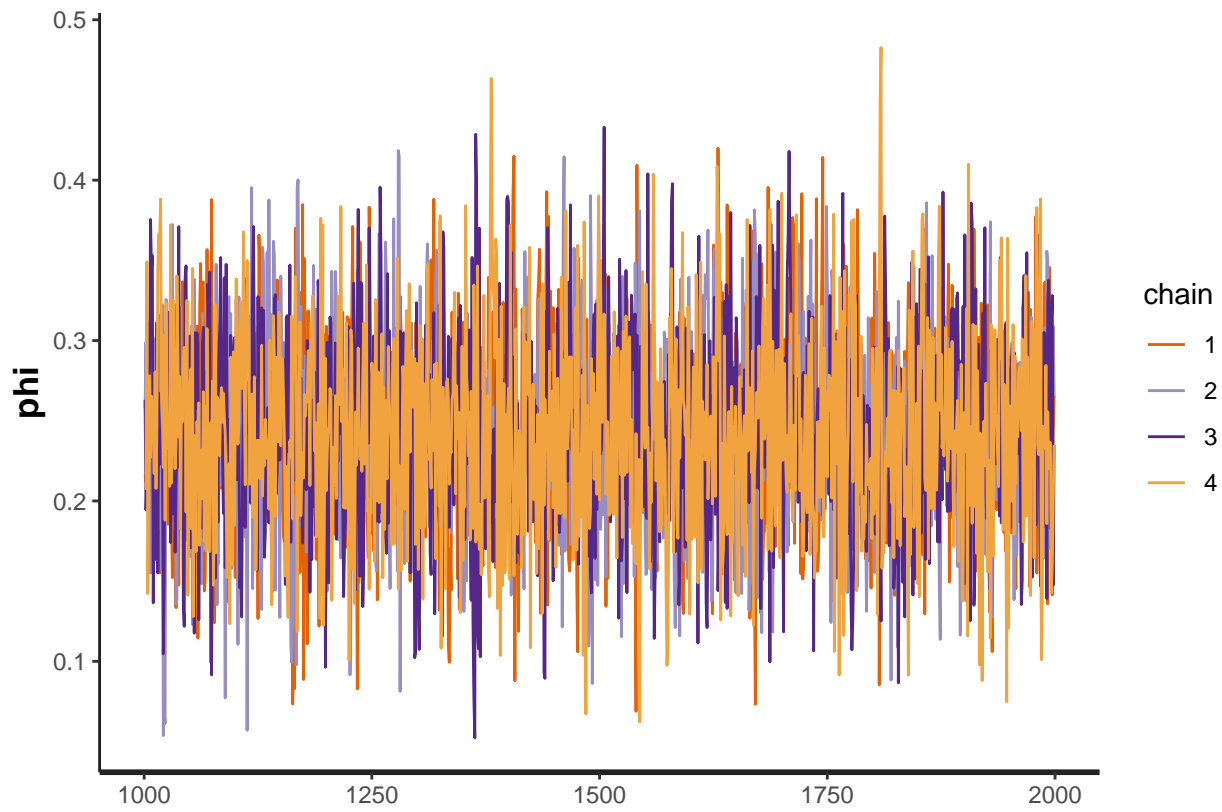
for the second dataset from $\phi=0.97$,

parameters	mean	95% credible intervals	the number of effective samples
mu	8.075	4.598-11.991	2518
phi	0.969	0.935-0.996	1881
sigma	1.914	1.749-2.093	2538

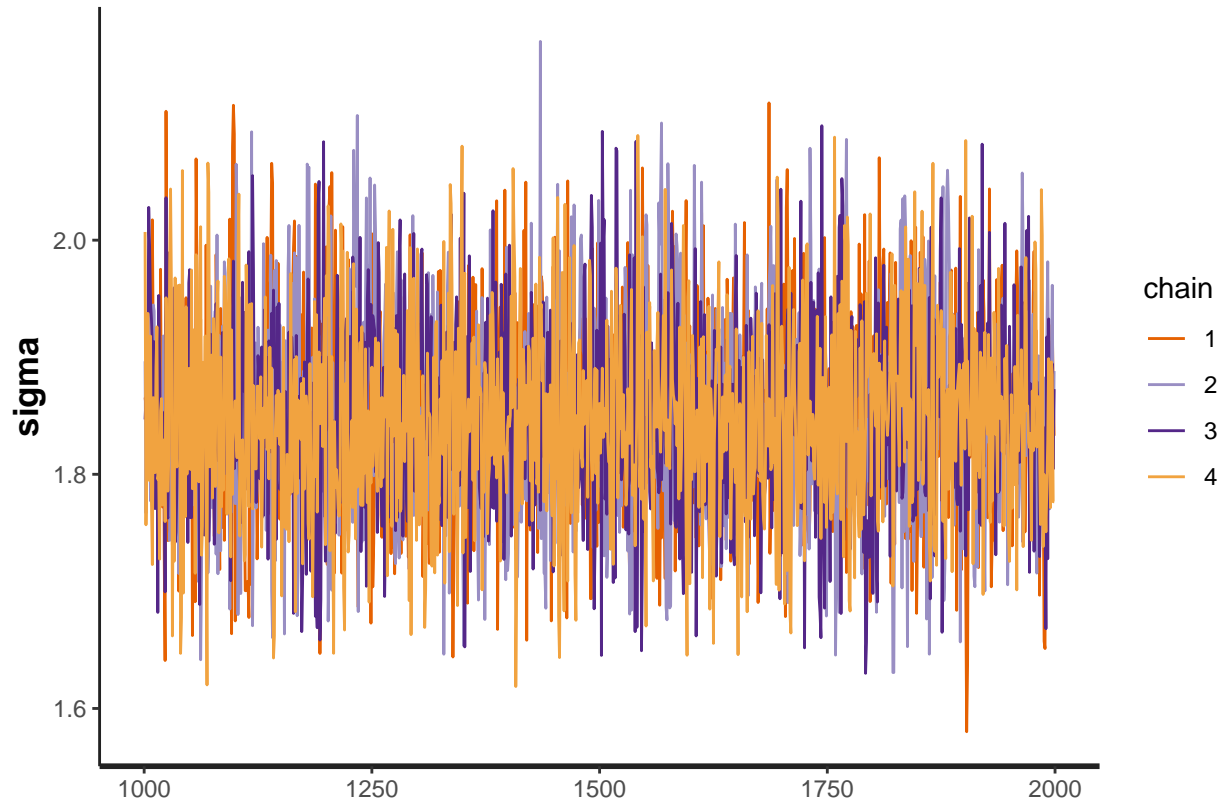
```
stan_trace(fit_x1, pars=c("mu"))
```



```
stan_trace(fit_x1, pars=c("phi"))
```



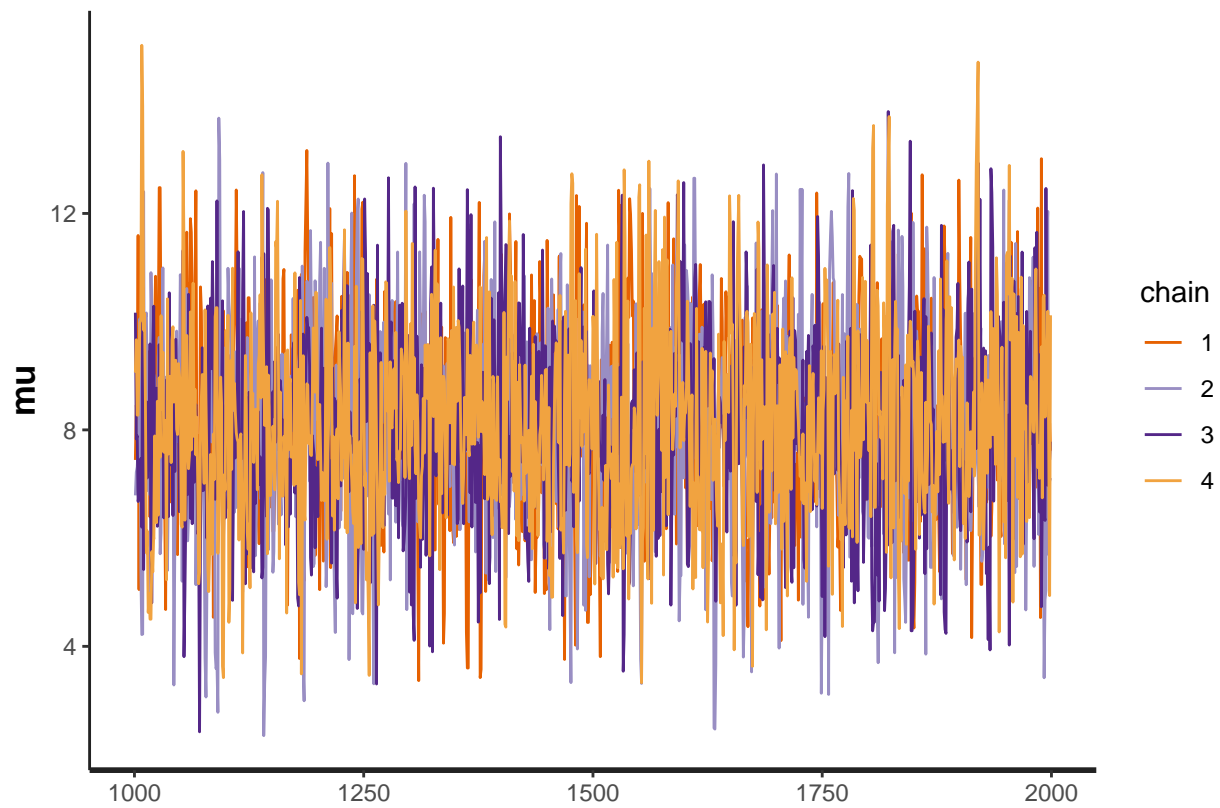
```
stan_trace(fit_x1, pars=c("sigma"))
```



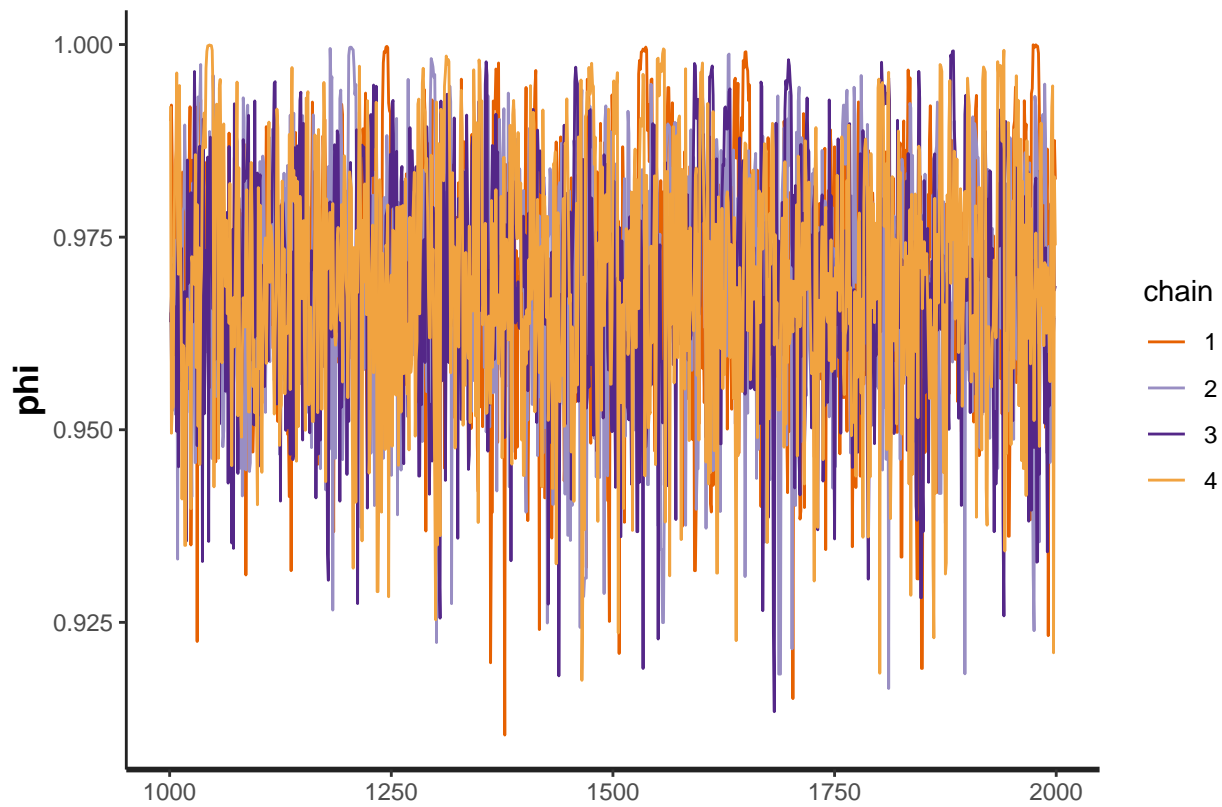
Above three trace plots are for the first dataset from $\phi=0.3$, it can be seen that all three parameters converge

quite well in the stationary distribution, the plots are also symmetric around the mean value.

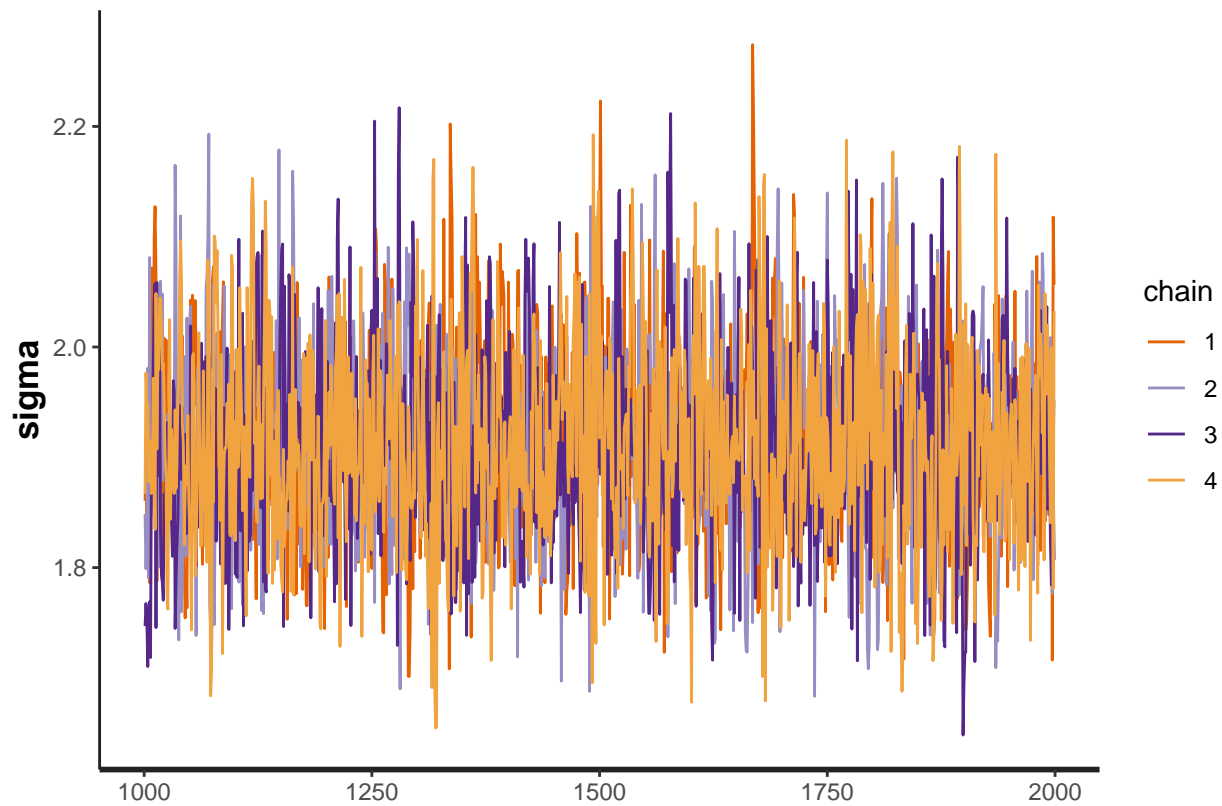
```
stan_trace(fit_y1, pars=c("mu"))
```



```
stan_trace(fit_y1, pars=c("phi"))
```



```
stan_trace(fit_y1, pars=c("sigma"))
```

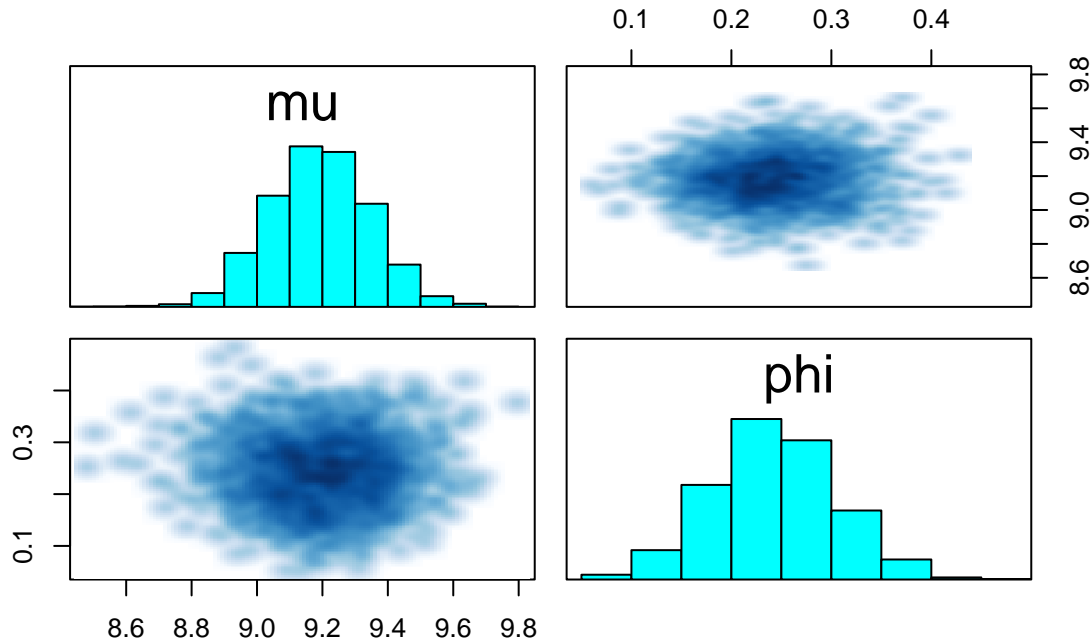


Above three trace plots are for the second dataset from $\phi=0.97$, it can be seen that the parameter μ

doesn't converge well, the variance of the points is bigger. This is because high ϕ leads to a strong positive autocorrelation.

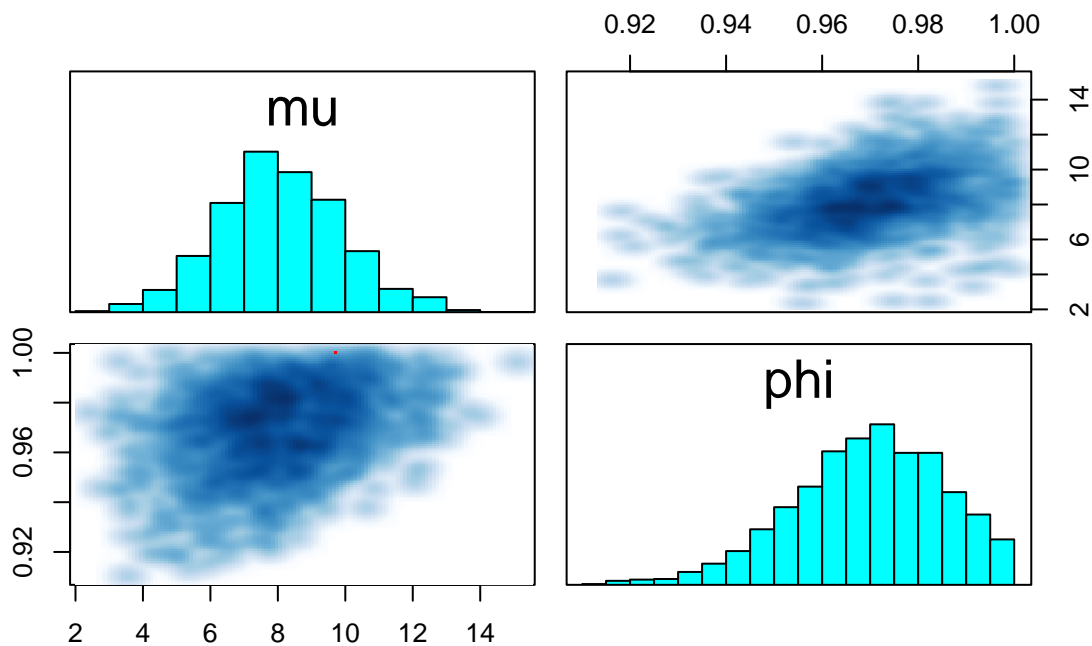
```
pairs(fit_x1, pars = c("mu", "phi"), main = "The first dataset from phi=0.3")
```

The first dataset from $\phi=0.3$



```
pairs(fit_y1, pars = c("mu", "phi"), main = "The second dataset from phi=0.97")
```

The second dataset from $\phi=0.97$



From above joint posterior plots of μ and ϕ , the first dataset generates symmetric ellipsoid shape, which could be possibly a multivariate normal distribution. While the second dataset generates a deformed ellipsoid,

which are hard to guess the distribution.