

Estimation scheme for fitting the stochastic model for the transmission of 2019-nCov in Hubei

Andrew Tredennick & John Drake

March 24, 2020

Introduction

This document describes a proposed estimation scheme for fitting the stochastic model (described elsewhere) of COVID-19 transmission to data. We plan to use the “synthetic likelihood” approach described by Wood (2010) to estimate unknown parameters. The posterior distribution of parameters will be explored using MCMC. The analysis presented here relies on the `synlik` R package functions. While the end goal is estimate parameters given real data, we start by using the model itself to simulate a trajectory of disease transmission and tuning the model fitting process to ensure we can recover known parameters.

Simulated data

Here we simulate data from the stochastic model using the `synlik` functions. The `synlik` functions are essentially wrappers around the core model functions written by Drake and Rohani. See the R script `simulators.R` for simulation wrappers that coerce the stochastic model into a form usable by `synlik`.

First, we define a `synlik` object using the pre-written wrappers.

```
# First source the necessary functions and load the synlik package
library(synlik)
source("stochastic-model.R") # the stochastic disease transmission model
source("simulators.R")      # simulator wrappers
source("style-definitions.R") # colors

# Create "synlik" object

# Parameters to estimate
params <- list(beta0 = 0.657, beta.factor = 1)

# Initial conditions, specific to Hubei
init <- list(S=59002000, E1=0, E2=0, E3=0, E4=0, E5=0, E6=0,
             I1 = 1, I2= 0, I3=0, I4=0, Iu1=0, Iu2=0, Iu3=0, Iu4=0,
             H=0, Ru=0, C=0)

# Time spans for simulation
start <- as.Date("2019-12-01")
today <- Sys.Date()

# Extra arguments needed for the stochastic model, specific to Hubei
extraArgs <- list("init" = init, "nstep" = NULL, "start" = start,
                  "today" = today, "dt" = 0.05, "w" = 40, "z" = 45, "c" = 1,
                  "presymptomatic" = 0, "timesToObs" = FALSE, "sigma" = 1/6.4,
                  "b" = 0.143, "a0" = 0.0446, "paramNames" = names(params), "nObs" = NULL)

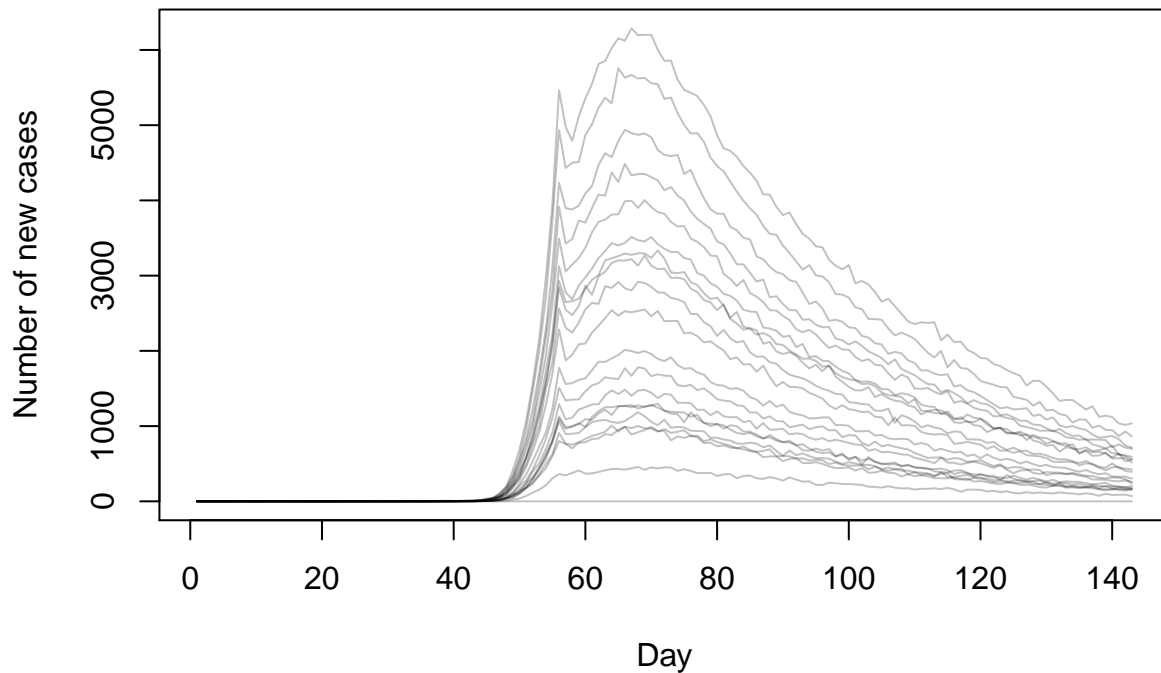
covid_sl <- new("synlik",
               simulator = mod_wrap,
```

```
param = log(unlist(params)),
extraArgs = extraArgs)
```

Then we simulate from the model to make sure it is working.

```
res <- simulate(covid_sl, nsim = 20, seed = 1234)
matplot(t(res), type = "l", col = col.cases.ci, lty = 1,
        main = "Simulated trajectories of COVID-19",
        xlab = "Day", ylab = "Number of new cases")
```

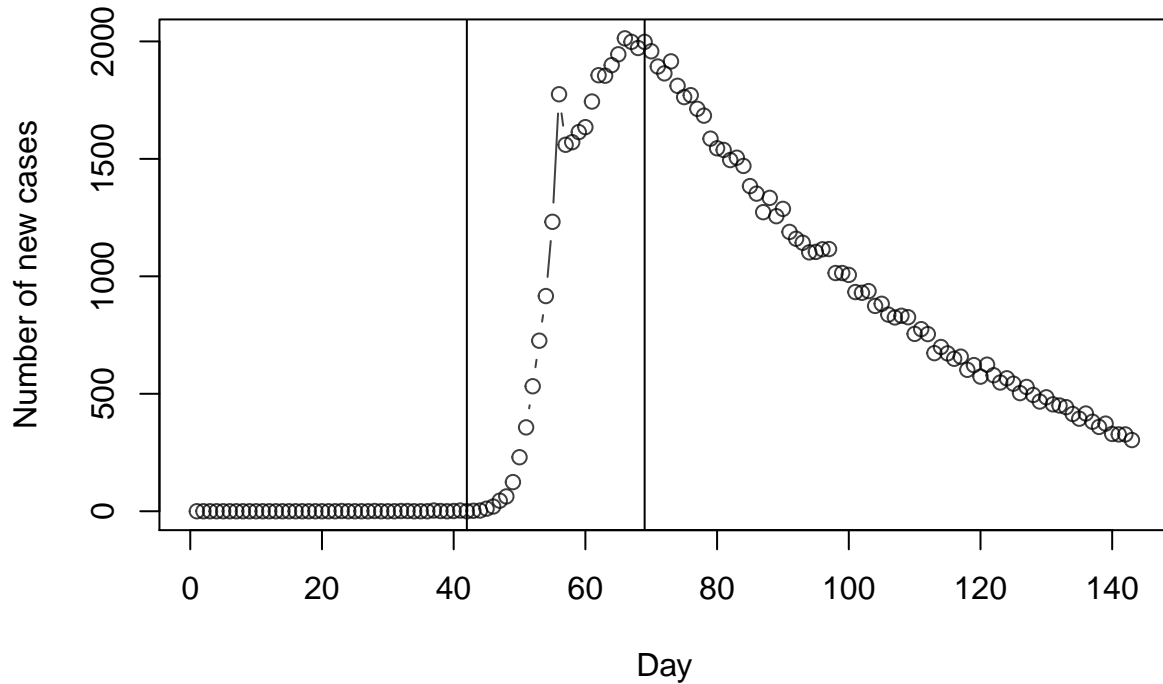
Simulated trajectories of COVID-19



Lastly, we can save one simulated trajectory as “data” for testing of the estimation scheme.

```
covid_sl@data <- res[1, ] # one replicate of the model simulations
covid_sl@extraArgs$obsData <- res[1, ]
covid_sl@extraArgs$nObs <- length(res[1, ])
plot(covid_sl@data, type='b', xlab='Day', col = col.cases,
     ylab='Number of new cases', main='Simulated COVID-19 cases in Hubei')
abline(v = c(42, 69))
```

Simulated COVID-19 cases in Hubei



Summary statistics for the synthetic likelihood

Fitting models using the synthetic likelihood (Wood 2010) requires calculating a vector of summary statistics (S) that summarize the observed time series and simulated trajectories. The summary statistics (S_{sim}) from simulated time series are compared to the summary statistics of the observed time series (S_{obs}) to evaluate the (synthetic) likelihood. For our particular problem, I used the following summary statistics based on phases of the epidemic delineated in the figure above by the solid vertical lines:

1. Maximum incidence before day 42 (stuttering chain of transmission)
2. Maximum incidence between day 42 and day 69 (exponential phase)
3. Maximum incidence after day 69 (decling phase)
4. Cumulative incidence before day 42
5. Cumulative incidence between day 42 and day 69
6. Cumulative incidence after day 69
7. Maximum incidence along the whole trajectory
8. Day of epidemic peak (timing of max incidence)
9. Final size of epidemic (total incidence)

Now we code up calculations for these summary statistics for `synlik`.

```
covid_stats <- function(x, extraArgs, ...) {  
  ## obsData is a vector of observed path  
  ## x is a M by n.t matrix of paths, each row of 'x' is a replicate  
  
  obsData <- extraArgs$obsData  
  
  stopifnot(is.vector(obsData), length(obsData) != 0)
```

```

if (!is.matrix(x)) x <- matrix(x, 1, length(x))

# Max incidence
X0 <- (apply(x, 1, function(x) max(x[1:41])))
X0 <- cbind(X0, (apply(x, 1, function(x) max(x[42:69]))))
X0 <- cbind(X0, (apply(x, 1, function(x) max(x[70:length(x)]))))

# Cumulative incidence
X0 <- cbind(X0, (apply(x, 1, function(x) sum(x[1:41]))))
X0 <- cbind(X0, (apply(x, 1, function(x) sum(x[42:69]))))
X0 <- cbind(X0, (apply(x, 1, function(x) sum(x[70:length(x)]))))

# Max along whole trajectory
X0 <- cbind(X0, apply(x, 1, max))

# Day of max
X0 <- cbind(X0, apply(x, 1, function(x) which.max(x)))

# Final epidemic size
X0 <- cbind(X0, apply(x, 1, sum))
return(X0)
}

```

And add them to the `synlik` object.

```
covid_sl@summaries <- covid_stats
```

Here are example statistics from a few simulations.

```
simulate(covid_sl, nsim = 5, stats = TRUE)
```

```
##      X0
## [1,]  2 1007  971  5 14755  34625 1007 64  49385
## [2,]  5 2918 2853 16 43659  97892 2918 68 141567
## [3,]  1  0  0  1  0  0  1  7  1
## [4,]  3 4939 4841 14 72684 168112 4939 66 240810
## [5,]  5 4004 3912 20 59526 140265 4004 69 199811

```

And the statistics from the “data”.

```
covid_sl@summaries(x = covid_sl@data, extraArgs = covid_sl@extraArgs)
```

```
##      X0
## [1,]  3 2013 1958 12 29695 69979 2013 66 99686

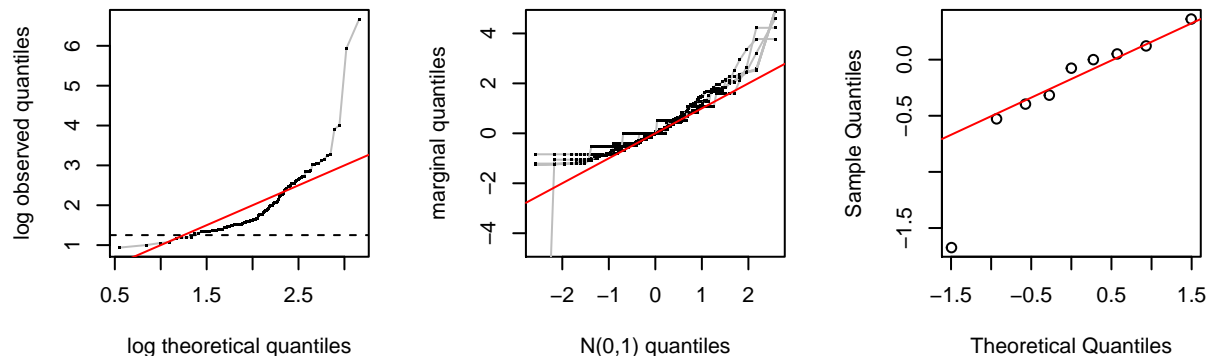
```

We also want to make sure the summary statistics meet the assumption of normality.

```
test_params <- c(beta0 = 0.43, beta.factor = 1)
checkNorm(covid_sl, param = log(test_params), nsim = 100)
```

```
##
## proportion |log(z)-log(q)|>.25 = 0.5

```



Model fitting

We are now ready to use MCMC to estimate unknown parameters. For the preliminary tests, we will try to fit two parameters: β_0 and ξ_β , which is the factor by which β_0 is reduced after policies designed to limit transmission were imposed on day x .

It is useful to make sure the likelihood can be evaluated before running MCMC.

```
test_params <- c(beta0 = 0.43, beta.factor = 1)
slik(covid_sl, param = unname(log(test_params)), nsim = 50)
```

```
## [1] -16393.77
```

We want to work with informative priors, so we define those for `synlik`. The prior for β_0 is $N(-0.4, 1)$; the prior for ξ_β is $N(-1.6, 0.2)$. Note that parameters are specified on the log scale for estimation and transformed to the arithmetic scale for simulation.

```
# This follows standard synlik notation
covid_priors_sl <- function(input, ...) { sum( input ) +
  dnorm (input[1], log(0.67), 1, log = TRUE) + # beta0 prior
  dnorm (input[2], log(0.2), 0.2, log = TRUE) } # beta reduction prior
```

Finally, we can run the `synlik::smcmc` routine.

```
# Define perturbation covariance matrix for MCMC proposals
covMat <- diag(rep(0.1, length(params)))

# Set the initial parameters, arithmetic scale
init_params <- c(beta0 = 0.5, beta.factor = 1)

# Set up a cluster for parallel simulations within MCMC steps
num_cores <- parallel::detectCores() - 1
cl <- parallel::makeCluster(num_cores)
parallel::clusterExport(cl, ls())
parallel::clusterEvalQ(cl, {
  suppressMessages(suppressWarnings(library(synlik)))
})
```

```
## [[1]]
## [1] "synlik"      "Rcpp"        "stats"       "graphics"    "grDevices"   "datasets"
## [7] "utils"      "methods"     "base"
```

```
##
## [[2]]
## [1] "synlik"      "Rcpp"      "stats"      "graphics"   "grDevices" "datasets"
## [7] "utils"      "methods"   "base"
##
## [[3]]
## [1] "synlik"      "Rcpp"      "stats"      "graphics"   "grDevices" "datasets"
## [7] "utils"      "methods"   "base"
##
## [[4]]
## [1] "synlik"      "Rcpp"      "stats"      "graphics"   "grDevices" "datasets"
## [7] "utils"      "methods"   "base"
##
## [[5]]
## [1] "synlik"      "Rcpp"      "stats"      "graphics"   "grDevices" "datasets"
## [7] "utils"      "methods"   "base"
##
## [[6]]
## [1] "synlik"      "Rcpp"      "stats"      "graphics"   "grDevices" "datasets"
## [7] "utils"      "methods"   "base"
##
## [[7]]
## [1] "synlik"      "Rcpp"      "stats"      "graphics"   "grDevices" "datasets"
## [7] "utils"      "methods"   "base"
```

Run the MCMC

```
covid_mcmc <- smcmc(covid_sl,
  initPar = log(init_params),
  nsim = 50,
  niter = 50,
  burn = 10,
  priorFun = covid_priors_sl,
  propCov = covMat,
  multicore = TRUE,
  ncores = num_cores,
  cluster = cl)
```

Stop the cluster

```
parallel::stopCluster(cl)
```

Plot the chains

```
par(mfrow = c(1, 2))
plot(covid_mcmc@chains[,1], type = "l", xlab = "Iteration",
  ylab = expression(log(beta)), las = 1)
plot(covid_mcmc@chains[,2], type = "l", xlab = "Iteration",
  ylab = expression(log(xi)), las = 1)
```

