

# Predicting the 2014 Ebola Outbreak in West Africa using Network Analysis

## Milestone Report

Shafi Bashar, Mike Percy, Romit Singhai

*{shafiab, mp81, romit}@stanford.edu*

### Abstract

The current Ebola outbreak in West Africa is the worst in history. Most traditional epidemiological models are compartmental models that have a random-mixing assumption. These models calculate the effective reproductive rate of an outbreak. We survey three of these models: the classic SIR (Susceptible, Infectious, Recovered) model and two extensions used in Ebola research.

Network models allow for avoiding the random-mixing assumption inherent in compartmental models. This is done by assigning each individual a finite set of permanent contacts. We review generated contact network models for SARS, including an urban network, a random network, and a scale-free network. We then review a worldwide network model that represents traffic flowing across transportation networks and consider approaches for predicting the extent of an Ebola outbreak.

### I. INTRODUCTION

TODO: Add introduction

### II. RELATED WORK

Related work by (Gomes et al., 2014) attempts to predict the spread of Ebola to different parts of the world based on a model that incorporates both the compartmental approach and the use of world-wide air traffic flows.

In (Meyers et al., 2005), the authors model the spread of the 2002-2003 outbreak of SARS in Hong Kong and Canada using a contact network. A contact network model attempts to characterize every interpersonal contact that can potentially lead to disease transmission in the community, with each person in the community represented as a node and each contact represented as an edge between them.

The majority of research in epidemiological theory is based on the compartmental model, which is not a network model. In order to simply capture the dynamics of disease spread over time, the compartmental model employs a population-wide random mixing assumption, meaning that each individual has a small and equal chance of coming into contact with any other individual in the population. To model the progress of an epidemic in a large population, the individuals in the population are compartmentalized according to the state of the disease. The most widely used such model is the SIR model introduced in (Kermack and McKendrick, 1932):

- **Susceptible (S):** Individuals who have not yet caught the disease from contact with an infectious individual.
- **Infectious (I):** Individuals who have the disease. They have some probability of infecting susceptible people.
- **Recovered (R):** Individuals who have experienced the full infectious period, and are now non-infectious and immune.

The changes among these states over time are represented by a set of differential equations. The basic reproductive number  $R_0$  is defined as the average number of secondary cases generated by a primary case in a pool of mostly susceptible individuals, and is an estimate of epidemic growth at the start of an outbreak if everyone is susceptible.

In (Chowell et al., 2004), the authors model the effect of Ebola outbreaks in 1995 in Congo and in 2000 in Uganda using a compartmental model similar to the SIR model. However, a distinct feature of Ebola is that individuals exposed to the virus who become infectious do so after a mean incubation period. In order to reflect this feature, the SIR model is extended with an additional “Exposed” compartment state. This SEIR model is summarized in Figure 1:



Fig. 1: SEIR model

In (Legrand et al., 2007), the Ebola outbreaks in Congo in 1995 and Uganda in 2000 are also studied. However, a major difference from (Chowell et al., 2004) is that (Legrand et al., 2007) models the spreading of disease in heterogeneous settings.

In order to gain better insight of the epidemic dynamics, the infectious phase is subdivided into three stages: infection in a community setting (I), infection in a hospital setting (H), and infection after death assuming a traditional funeral (F).

This SEIHFR compartmental model is summarized in Figure 2.

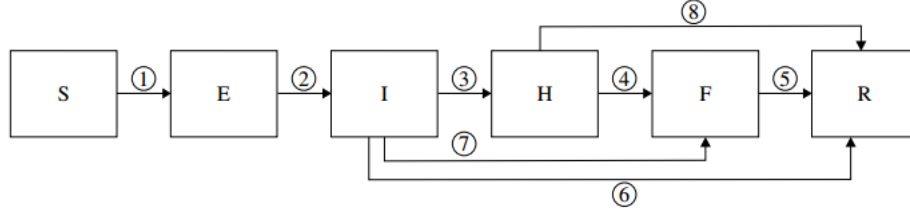


Fig. 2: SEIHFR model

For our purposes, as described below, we stopped short of implementing this complicated model.

### III. MODELING LOCALIZED EPIDEMIC SPREAD

#### *Estimating model parameters and basic reproduction number for 2014 Ebola data using a random mixing model*

In first phase of the project, to calculate the basic reproduction number of the current Ebola epidemic spread at different countries, we performed model fitting on the data we gathered from Rivers (2014). We considered the SEIR model described in Chowell et al. (2004). The SEIR model under consideration is a non-linear model with six parameters. The current Ebola epidemic is still spreading and depending on the preventative measures taken, the underlying dynamics of the spread can change drastically anytime. We fitted our model to three countries in the West Africa - Guinea, Sierra Leone and Liberia. In addition, we performed model fitting for the West Africa Region by adding up the data from these three countries. Given the limited number of data available, instead of fitting all six parameters to the model, we decided to fix some of the parameters based on the studies on previous Ebola epidemic. In Chowell et al. (2004), the incubation time of the Ebola  $1/k$  is found to be varying between 1 to 21 days, with a mean time of 6.3 days for previous Ebola spread. For ease of data fitting, we set this parameter value to the mean value of 6.3 days. We note that, the dynamics of the current epidemic may differ from previous one, and therefore fixing a value based on the prior estimate may lead to some inaccuracy. To model the effect the intervention on the spread of the epidemic, a modified transmission rate  $\beta_1$  is generally used instead of the initial transmission rate  $\beta_0$  at later stage of the epidemic. The transition of  $\beta_0$  to  $\beta_1$  depends on intervention time  $\tau$  and decaying factor  $q$ . The choice of intervention time is a very difficult problem. To this end, we looked into different sources like Wikipedia, WHO and CDC website to learn more about the timeline of the spread. In Guinea, a 2-year-old boy fell in December 2, 2013, later diagnosed as Ebola patient. We consider this incidence as the index case for Guinea and set  $t_0$  to December 2. In March 2, the Government of Guinea informed WHO regarding the possibility of Ebola epidemic and declared national health emergency. We considered this date as the intervention date and set  $\tau$  to 110. In Sierra Leone, one person fell in April 2014. In June 12, 2014 the country declared emergency and closed borders with neighboring Guinea and Liberia. We consider the first date as  $t_0$  and second date as intervention time, therefore set  $\tau$  to 50. In Liberia, in March 31, 2014, there were official confirmation of two person getting infected from Ebola. We set this date to  $t_0$  and set  $I_0$  and  $C_0$  in our model as 2. The Government of Liberia shut down all schools in July 30, 2014. We consider this date as the date of intervention and set  $\tau$  to 120.

In order to fit the non-linear SEIR model to data, we used non-linear least square estimation. We use the reported data  $(t_i, c_i)$  for  $i = 1, 2, \dots, n$  where  $t_i$  denote the  $i$ -th reporting time and  $c_i$  as the cumulative number of infectious cases from the beginning of outbreak time  $t_0$  to time  $t_i$ . The optimization problem contains a large number of local minimas. Therefore, the choice of initial parameter estimate is an important consideration to get the global optimum solution. In order to find a good initial choice of the parameters as an input to the non-linear least square solver, we first perform a Latin hypercube sampling on the 4-dimensional parameter space. We grid up the hypercube with a number of grid points in each dimension. We then choose the sample that minimizes the least square error as the initial input. In order to calculate the 95% confidence interval of estimated parameter, we performed bootstrapping based on residual error.

The estimated parameters for the SEIR model of Guinea, Sierra Leone, Liberia, as well as West Africa region is presented in the appendix. In Figure 3, we presented the the number of incidence at different compartments of the SEIR model, as well as the cumulative number of infectious cases over time. Our model fit was based on the last data we collected in October 21, 2014. As of today, November 11, 2014, few additional data points are available; we also plotted those data points in the graph. In addition to that, we extrapolated the graph up to December 31, 2014. We note that, the forecasting to future cases may not be appropriate as the underlying factors of the Epidemic are changing rapidly with the increase in safety measures. We observe that, the prediction of the model to the cases up to November 11 is mostly in par with the observed data for Guinea and Sierra Leone. We note that, for Guinea epidemic, we have data for the longest range of days as the initial case was in December 2. The estimated model, therefore, captures most of the dynamics of the spread in Guinea in contrast to

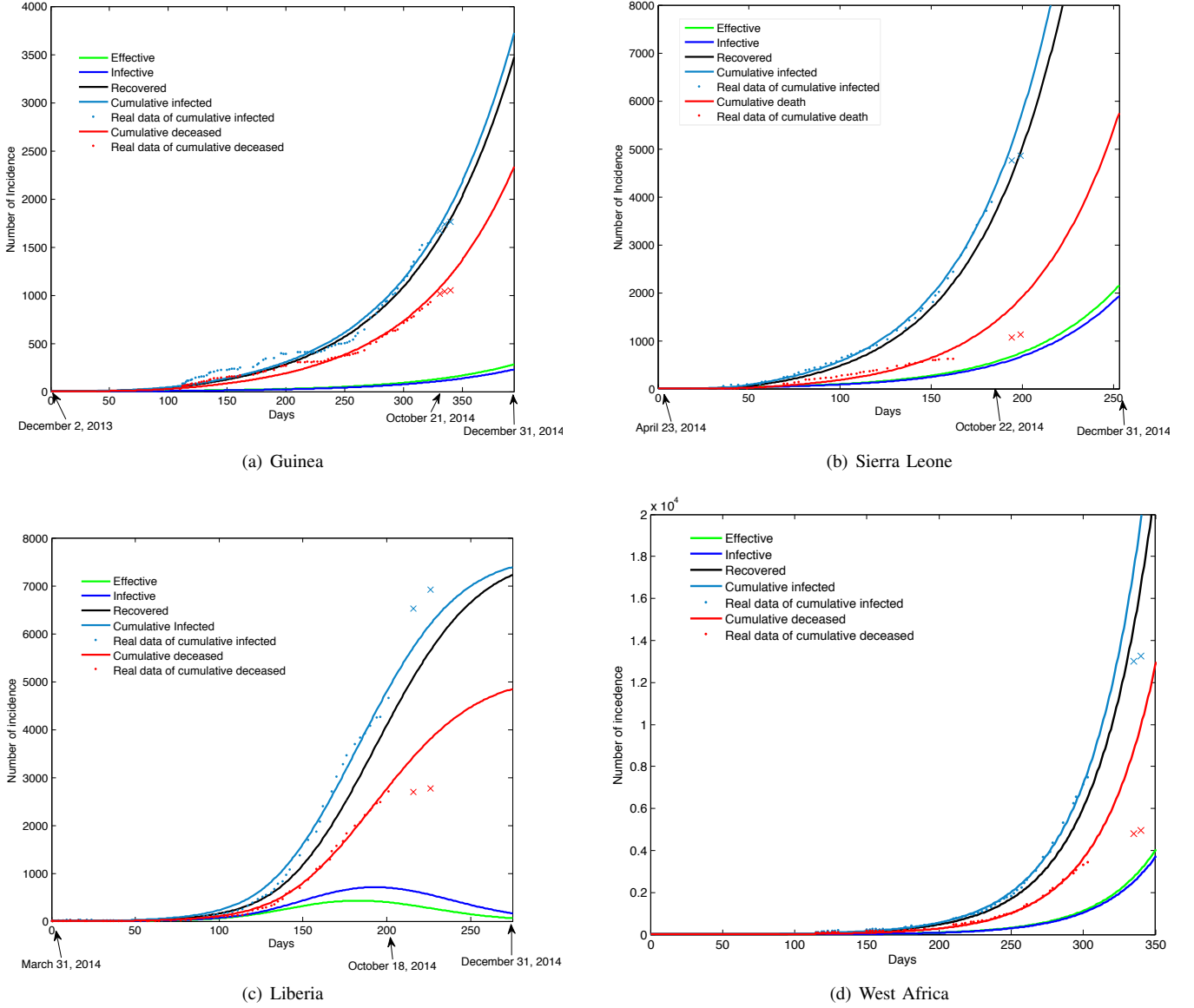


Fig. 3: SEIR model fit results for 2014 Ebola epidemic data

the other cases. In case of Liberia, our prediction under-estimated the observed data. Our guess is the estimated parameters for Liberia may be overfitting the model. While performing some exploratory analysis on the Liberia data, we observed some discrepancies, e.g. decrease in cumulative value from the previous data points, which may somewhat distorted the model fitting parameters. For the West Africa region plot, our prediction number is much higher than the actual data. This is understandable, as the underlying assumption behind SEIR model is random mixing. Within a country without any movement restriction, this model is somewhat appropriate. However, for the aggregated model among multiple countries, this assumption is no longer valid due to stricter movement regulations between country borders. Therefore a model with random mixing assumption will over predict the infectious number of cases.

As an extension of the above approach, to find parameters for the model based on the current data, we designed and ran the simulation using mcmc approach with input as time series data for the various infected regions. Initial values for the various parameters were assumed as above for the prior distribution but the range based on the literature was used to find the most appropriate values for the current outbreak. For example  $\beta_0$  had a range  $[0,1]$ ,  $\beta_1$  had a range  $[0,1]$ , infection time ( $1/\gamma$ ) had range  $[3.5,10.7]$ , incubation time ( $1/k$ ) had a range  $[5,22]$  and ( $\tau$ ) had a range  $[100,150]$ . The simulation is currently running on a server and we will have results in the final report.

### Mapping the compartmental model to the network based model

It is preferable to use a model that does not assume random mixing. One such network model we are integrating is a data set from a social networking site from the [UCI Machine Learning Repository \(2008\)](#). This undirected network has 4,846,609 nodes and 42,851,237 edges and average degree of 17.7. It is clearly a scale-free network because the degree distribution follows a power law after a degree of approximately 50. Using this data set, we calculate transmissibility and epidemic threshold.

We calculated the transmissibility  $T$  of a disease which is defined as the average probability that an infectious individual will transmit the disease to a susceptible individual with whom they have contact. We also calculated the epidemic threshold  $T_c$  which is the minimum transmissibility required for an outbreak to become a large-scale epidemic.

We start with the reproductive number  $R_0$  calculated using the SEIR compartmental model. Using the relationship between the basic reproductive number  $R_0$  and the transmissibility  $T$  given in [Meyers et al. \(2005\)](#):

$$R_0 = T \frac{\langle k^2 \rangle}{\langle k \rangle - 1}$$

while using values  $R_0 = 2.27$  and  $k = 17.7$  we find that transmissibility  $T = 0.1210$  for this network.

We also calculate the value of epidemic threshold  $T_c$ , which is defined by Meyers as:

$$T_c = \frac{\langle k \rangle}{\langle k^2 \rangle - \langle k \rangle}$$

for mean degree  $\langle k \rangle$  and mean square degree  $\langle k^2 \rangle$ . Our calculations show that  $T_c = 0.0598$  for this network.

We are using these results to predict average size of the outbreak and to compare the behavior of a random mixing compartmental model and a network model. We are incorporating this as a local model into our worldwide simulations.

## IV. PREDICTING WORLDWIDE SPREAD

### Approach

In order to simulate how the Ebola epidemic might spread across the world, we have assembled a worldwide network based on international trade. Because the trade numbers are in US dollars, we have assumed a linear relationship between exports in dollars and travelers going abroad. This is a strong assumption, and we are considering more sophisticated mappings. We are using this as a long-range network representing trade-based population movement and connecting localized subnetworks with these trade-driven edges. We are developing a simulation framework to run in discrete time steps predict the spread of Ebola across this worldwide network.

### Trade network data preparation

The trade network dataset we are using is based on data retrieved from the web service API to the UN Comtrade database ([United Nations, 2014](#)) via their web service API. The data queried for were country-to-country SITC-1 exports from the latest data available for each country.

Several problems with this raw dataset immediately became apparent which required working around. Many war-torn countries such as Liberia have not reported detailed export data to the UN in decades (Liberia last reported in 1984). As a result of this, the export totals are incorrect for the present day. In addition, this old data reports exports to countries that no longer exist, including East Germany and Yugoslavia. In order to make the dataset usable for modern-day predictions, we performed the following transformations on the data (using Perl scripts):

- Manually remapped non-existent countries to their modern equivalents. In some cases, we had to do a best-effort mapping. For example, Yugoslavia split into many states, so we remapped exports to Yugoslavia be exports to modern-day Slovakia, since Slovakia currently has the largest economy of the states that once comprised Yugoslavia.
- Removed exports from states that no longer exist, preferring export data from modern equivalents.

We then needed to make the per-country exports sum up to the latest available export data. The United States [Central Intelligence Agency \(2013\)](#) World Factbook contains up-to-date export totals for most countries in the world. Using the above data set, on the (admittedly strong) assumption that the export distribution from country-to-country in the UN Comtrade data set has remained the same between the last year a country has reported detailed export data and the latest totals from the CIA, we linearly renormalized each outgoing edge in our data set so that the total sum of exports equalled the latest CIA data for each country. This required a manual step of mapping country names in the UN dataset to those used in the CIA dataset.

Once we had a “renormalized” dataset including detailed export edges and totals matching the latest available data, we needed to map these numbers to theoretical international travelers. We could not find complete, or even nearly complete, publicly available data on the number of international outgoing travelers per country. While inbound tourism numbers are available from [The World Bank \(2014b\)](#), for the outbound tourism numbers from [The World Bank \(2014c\)](#), many countries are missing, especially the African nations that we care so much about from an Ebola outbreak perspective. Some data are available from the [UN World Tourism Organization \(2014\)](#), however they appear to be behind a paywall. Data for total air

travellers are available from [The World Bank \(2014a\)](#), however this total includes domestic flights and so is less useful for our purposes.

Our current approach to map outbound dollars to outbound travelers is to assume a linear relationship between these numbers, and also to assume that the same linear relationship holds for imports to inbound tourists to a country. We currently use the United States as a model and use the ratio of imports to the United States per year to the number of tourists visiting the United States per year. Based on data from the [US Office of Travel & Tourism Industries \(2013\)](#), 69.77 million people visited the United States in 2013. According to our renormalized data set, total imports into the United States were \$2.21 trillion during the same period. Dividing imports by visitors (an admittedly simplistic approach) gives us a scaling factor of approx. 31,665. Therefore we have applied this scaling factor to all edges in our international exports network, giving us some approximation of outbound travellers from country to country based on export numbers.

### ***Supplementing the trade network***

Using a trade network has its flaws, even discounting errors stemming from poor approximations. For one, a trade-centric network ignores activity with low economic impact, such as a vegetable farmer traveling to a nearby country to sell his products at the market, or someone driving across the border to visit a family member. Relative to the economic impact of the industrial diamond or rubber trade (exports of Liberia) for example, these potentially disease-spreading behaviors simply are not represented equitably if at all in the model.

On the other hand, economic trade data is widely accessible and is likely to be fairly accurate, and certainly it is safe to assume that that trade activity correlates with travel between two countries.

Due to the above drawbacks, we plan to supplement our trade network with a generated network that captures local commute-related movement as well. One approach to generating such a network is the gravity law, noted in a survey by [Barthélemy \(2011\)](#), which claims that the number of trips from point A to point B is inversely proportional to the square of the distance between them. Such a model was used to study traffic on Korean highways in [Jung et al. \(2008\)](#).

### ***Simulation model***

We are currently developing a simulation framework to incorporate intra-country spreading behavior (using compartmental and network-based models) with inter-country links based on the above worldwide network models. A related approach is the GLEaM stochastic simulation model developed by [Balcan et al. \(2010\)](#). Our simulation model is a stochastic, discrete time-step model in which each time step is one day, and for each day, each infectious individual may travel to another country with some probability related to the worldwide networks. Within each country, infectious individuals spread their disease with some probability related to the intra-country model being used (compartmental or network-based).

## **REFERENCES**

- Duygu Balcan, Bruno Gonçalves, Hao Hu, José J Ramasco, Vittoria Colizza, and Alessandro Vespignani. Modeling the spatial spread of infectious diseases: The global epidemic and mobility computational model. *Journal of computational science*, 1 (3):132–145, 2010.
- Marc Barthélemy. Spatial networks. *Physics Reports*, 499(1):1–101, 2011.
- Central Intelligence Agency. The World Factbook 2013-14: Country Comparison: Exports, 2013. URL <https://www.cia.gov/library/publications/the-world-factbook/rankorder/2078rank.html>. [Online; accessed 6-November-2014].
- Gerardo Chowell, Nick W Hengartner, Carlos Castillo-Chavez, Paul W Fenimore, and JM Hyman. The basic reproductive number of Ebola and the effects of public health measures: the cases of Congo and Uganda. *Journal of Theoretical Biology*, 229(1):119–126, 2004.
- MF Gomes, AP Piontti, Luca Rossi, Dennis Chao, Ira Longini, M Elizabeth Halloran, and Alessandro Vespignani. Assessing the international spreading risk associated with the 2014 West African Ebola outbreak. *PLoS Currents Outbreaks*, 2014.
- Woo-Sung Jung, Fengzhong Wang, and H Eugene Stanley. Gravity model in the korean highway. *EPL (Europhysics Letters)*, 81(4):48005, 2008.
- William O Kermack and Anderson G McKendrick. Contributions to the mathematical theory of epidemics. II. The problem of endemicity. *Proceedings of the Royal society of London. Series A*, 138(834):55–83, 1932.
- J Legrand, RF Grais, PY Boelle, AJ Valleron, and A Flahault. Understanding the dynamics of Ebola epidemics. *Epidemiology and infection*, 135(04):610–621, 2007.
- Lauren Ancel Meyers, Babak Pourbohloul, Mark EJ Newman, Danuta M Skowronski, and Robert C Brunham. Network theory and SARS: predicting outbreak diversity. *Journal of theoretical biology*, 232(1):71–81, 2005.
- Caitlin Rivers. Data for the 2014 Ebola outbreak in West Africa, 2014. URL <https://github.com/cmrsivers/ebola>. [Online; accessed 15-October-2014].
- The World Bank. Air transport, passengers carried, 2014a. URL <http://data.worldbank.org/indicator/IS.AIR.PSGR>. [Online; accessed 12-November-2014].

- The World Bank. International tourism, number of arrivals, 2014b. URL <http://data.worldbank.org/indicator/ST.INT.ARVL>. [Online; accessed 11-November-2014].
- The World Bank. International tourism, number of departures, 2014c. URL <http://data.worldbank.org/indicator/ST.INT.DPRT>. [Online; accessed 12-November-2014].
- UCI Machine Learning Repository. Social Network from TopCoder's Epidemic Contest, 2008. URL <http://networkdata.ics.uci.edu/data.php?id=108>. [Online; accessed 12-November-2014].
- United Nations. United Nations Commodity Trade Statistics Database, 2014. URL <http://comtrade.un.org/data/>. [Online; accessed 6-November-2014].
- UN World Tourism Organization. Outbound tourism data (calculated on the basis of arrivals data in destination countries), 2014. URL <http://statistics.unwto.org/content/outbound-tourism-data-calculated-basis-arrivals-data-destination-countries>. [Online; accessed 12-November-2014].
- US Office of Travel & Tourism Industries. 2013 Monthly Tourism Statistics, 2013. URL <http://travel.trade.gov/view/m-2013-I-001/table1.html>. [Online; accessed 6-November-2014].

## APPENDIX

TABLE I: Parameter estimation for Ebola SEIR model (Guinea & Sierra Leone)

Incidence Dependent Parameters	Guinea			Sierra Leone		
	Value	Comments		Value	Comments	
Initial Case $t_0$	December 2, 2013	one person fell ill		April 23, 2014	one person fell ill	
$S_0$	0	-		0	-	
$E_0$	0	-		0	-	
$I_0$	1	-		1	-	
$R_0$	0	-		0	-	
$C_0$	1	-		1	-	
Intervention time	March 2, 2014	Gov. of Guinea informed WHO		June 12, 2014	Country declared emergency	
$\tau$	110	-		50	-	
Estimated Parameters	Value	95% CI	Comments	Value	95% CI	Comments
Incubation Time $1/k$	6.3	-	based on previous works	6.3	-	based on previous works
Infection Time $1/\gamma$	5.4957	[5.43, 5.545]	-	6.36	[6.324, 6.396]	-
$\beta_0$	0.2407	[0.2374, 0.244]	-	0.357	[0.3525, 0.3613]	-
$\beta_1$	0.2084	[0.2033, 0.2135]	-	0.2012	[0.1994, 0.2028]	-
$q$	32	[0.1, 100]	-	34	[1.9, 110]	-
Fatality Rate	0.67	-	-	0.38	-	-
$R_0$	1.323	[1.295, 1.341]	-	2.27	[2.243, 2.298]	-
$R_1$	1.145	-	-	1.28	-	-

TABLE II: Parameter estimation for Ebola SEIR model (Liberia & West Africa overall)

Incidence Dependent Parameters	Liberia			West Africa		
	Value	Comments		Value	Comments	
Initial Case $t_0$	March 31, 2014	official confirmation two infected		December 2, 2013	one person fell ill in Guinea	
$S_0$	0	-		0	-	
$E_0$	0	-		0	-	
$I_0$	2	-		1	-	
$R_0$	0	-		0	-	
$C_0$	2	-		1	-	
Intervention time	July 30, 2014	School shutdown		March 2, 2014	Gov. of Guinea informed WHO	
$\tau$	120	-		110	-	
Estimated Parameters	Value	95% CI	Comments	Value	95% CI	Comments
Incubation Time $1/k$	6.3	-	based on previous works	6.3	-	based on previous works
Infection Time $1/\gamma$	10.7	[8.9, 10.7]	-	6.8	-	-
$\beta_0$	0.169	[0.168, 0.191]	-	0.2	-	-
$\beta_1$	0.0001	[0.0001, 0.1]	-	0	-	-
$q$	0.0085	[0.0072-0.0264]	-	0	-	-
Fatality Rate	0.67	-	-	-	-	-
$R_0$	1.808	-	-	1.36	-	-
$R_1$	0.001	-	-	-	-	-