

# Modeling Epidemics: Introduction, Simple Model, and Linear Least Squares

Brian Hunt  
University of Maryland  
AMSC/MATH 420, Spring 2013

# First Models

- Preliminary goal: Model the spread of an infectious (contagious) illness through a population.
- Simplifying assumptions:
  - The total population  $N$  is constant in time.
  - A newly infected person becomes infectious the next day and remains infectious forever.
  - Each infectious person is equally likely (probability  $p$ ) to infect each noninfectious person on a given day.
- Let  $I(t)$  be the number of infectious people at the start of day  $t$ .

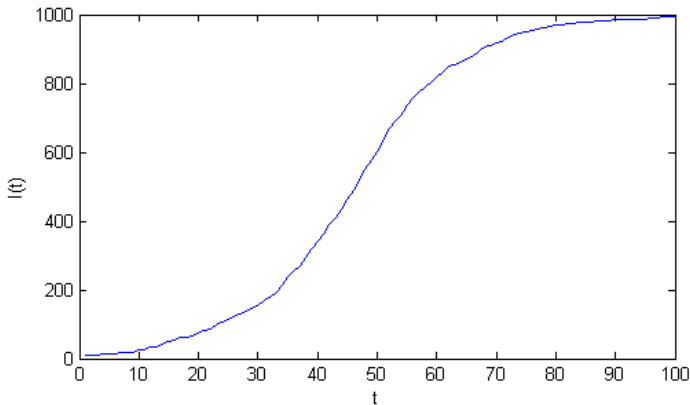
# Stochastic Model

- Number the people from 1 to  $N$ .
- Let  $x_n(t)$  be the infectious status (1 if infectious, 0 if not) of person  $n$  at the start of day  $t$ .
- We can simulate a possible spread of the illness with the following program ("rand"= random number):

```
for t=1:T-1
  for n=1:N
    let x(n,t+1)=x(n,t)
    for m=1:N
      if x(m,t)=1 and rand<p, then let x(n,t+1)=1
    end
  end
end
end
```

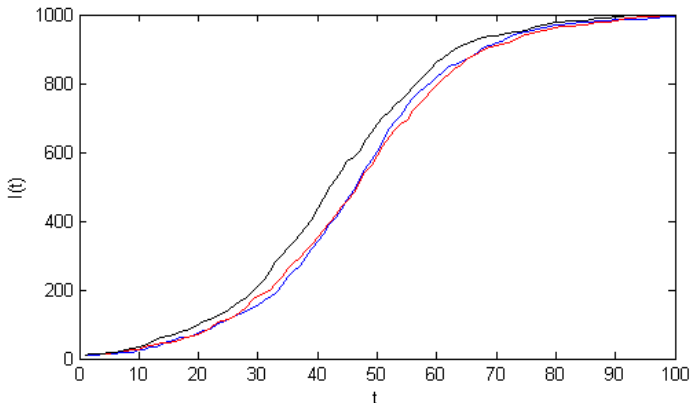
# Simulation Results

- Notice that  $\mathcal{I}(t) = \sum_{n=1}^N x_n(t)$ .
- Here are the results of a simulation with  $p = 10^{-4}$ ,  $N = 1000$ , and  $\mathcal{I}(1) = 10$ :



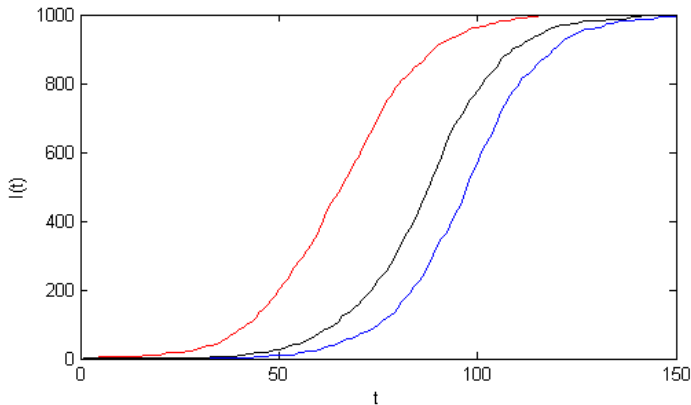
# Simulation Results

- And here are the results of three different simulations with  $p = 10^{-4}$ ,  $N = 1000$ , and  $\mathcal{I}(1) = 10$ :



# Simulation Results

- Finally, here are the results of three different simulations with  $p = 10^{-4}$ ,  $N = 1000$ , and  $\mathcal{I}(1) = 1$ :



## Expected (Average) Daily Outcome

- Let's determine the expected number of people infected on a day  $t$  that starts with  $\mathcal{I}(t)$  infectious people and  $N - \mathcal{I}(t)$  who are **susceptible** to infection.
- A susceptible person  $n$  has probability  $1 - p$  of NOT being infected on day  $t$  by a given infectious person  $m$ . Therefore, person  $n$  has probability  $(1 - p)^{\mathcal{I}(t)}$  of NOT being infected on day  $t$ .
- The expected number of people who are infected on day  $t$  is then  $[1 - (1 - p)^{\mathcal{I}(t)}][N - \mathcal{I}(t)]$ , so

$$E[\mathcal{I}(t + 1)] = \mathcal{I}(t) + [1 - (1 - p)^{\mathcal{I}(t)}][N - \mathcal{I}(t)]$$

# Deterministic Models

- If both  $\mathcal{I}(t)$  and  $N - \mathcal{I}(t)$  are large enough, it may be reasonable to approximate  $\mathcal{I}(t + 1)$  by its expected value, resulting in a deterministic model:

$$\mathcal{I}(t + 1) = \mathcal{I}(t) + [1 - (1 - p)^{\mathcal{I}(t)}][N - \mathcal{I}(t)] \quad (1)$$

- If  $p\mathcal{I}(t)$  is small, we can approximate  $(1 - p)^{\mathcal{I}(t)}$  by  $1 - p\mathcal{I}(t)$ , yielding a simpler model:

$$\mathcal{I}(t + 1) = \mathcal{I}(t) + p\mathcal{I}(t)[N - \mathcal{I}(t)] \quad (2)$$

- For these models, given  $\mathcal{I}(1)$  we can compute  $\mathcal{I}(2)$ ,  $\mathcal{I}(3)$ , ....



# Deterministic versus Stochastic

- These deterministic models are much more efficient to compute (1 calculation versus  $N^2$  for the stochastic model). Their predictions may be just as reasonable as any particular simulation of the stochastic model.
- The stochastic model can give some idea of the uncertainty of its predictions via multiple simulations; the deterministic models we've written down say nothing about their uncertainty.

# Continuous-Time Model

- The models we have discussed so far are called **discrete-time** models; time  $t$  takes on only integer values.
- We can approximate these models by continuous-time processes; approximating model (2), we get

$$\mathcal{I}'(t) = p\mathcal{I}(t)[N - \mathcal{I}(t)] \quad (3)$$

- We can write down an exact solution to this differential equation:

$$\mathcal{I}(t) = \frac{N\mathcal{I}(0)}{\mathcal{I}(0) + [N - \mathcal{I}(0)]e^{-pNt}}$$

# Fitting the Model to Data

- The solution  $\mathcal{I}(t)$  of model (3) has three parameters:  $N$ ,  $\rho$ , and  $\mathcal{I}(0)$ . Suppose we know  $N$  but not the other two parameters. Given a set of data points  $[t_j, \mathcal{I}_j]$ , we can ask which values of  $\rho$  and  $\mathcal{I}(0)$  best fit the data.
- [A more fundamental (but more difficult) question is whether the model can adequately fit the data at all; are there ANY parameters of the model that fit the data reasonably well?]
- We could try to minimize the sum of the squares of the residuals  $\mathcal{I}_j - \mathcal{I}(t_j)$ . However, this would be a NONlinear least squares problem, because  $\mathcal{I}(t)$  does not depend linearly on  $\rho$  or  $\mathcal{I}(0)$ .

# Way 1 to use Linear Least Squares

- If the data is given at consecutive values of  $t$ , say  $t_j = j$ , then one approach is to use model (2) and write

$$\mathcal{I}(t+1) - \mathcal{I}(t) = \rho \mathcal{I}(t)[N - \mathcal{I}(t)].$$

The right-hand side is a linear function of the parameter  $\rho$ , and linear least squares yields the value of  $\rho$  that minimizes the sum of the squares of the residuals  $\mathcal{I}_{j+1} - \mathcal{I}_j - \rho \mathcal{I}_j(N - \mathcal{I}_j)$ .

- This doesn't resolve the question of which value of  $\mathcal{I}(0)$  to use. If we let  $t_0 = 0$  for the first data point, then we could let  $\mathcal{I}(0) = \mathcal{I}_0$ . However, this might not be the best choice of  $\mathcal{I}(0)$  in order to make the residuals  $\mathcal{I}_j - \mathcal{I}(t_j)$  small.

## Way 2 to use Linear Least Squares

- Going back to the solution of model (3), we can make a transformation of variables so that the transformed solution does depend linearly on its parameters. First we divide both sides into  $N$  and simplify:

$$N/I(t) = 1 + [N/I(0) - 1]e^{-pNt}$$

- Next subtract 1 and take the logarithm:

$$\log[N/I(t) - 1] = \log[N/I(0) - 1] - pNt$$

- Let  $Z(t) = \log[N/I(t) - 1]$ ; then the model becomes  $Z(t) = Z(0) - pNt$ . This is a linear function of the parameters  $pN$  and  $Z(0)$ . One can transform the data to pairs  $(t_j, Z_j)$ , use linear least squares to determine values for  $pN$  and  $Z(0)$ , and then solve for  $p$  and  $I(0)$ .

# Caveat

- Both ways of using linear least squares transform the model or its solution into a linear relationship between two quantities that can be computed from the data points  $(t_j, \mathcal{I}_j)$ ; in the second way, the model predicts that  $Z_j$  is a linear function of  $t_j$ .
- Rather than simply accept the result of the least squares fit, one should graph the predicted relationship (e.g.,  $Z_j$  versus  $t_j$ ) and see if it actually looks linear. This gives some idea of how valid the model is.
- Regardless of how one determines values for  $p$  and  $\mathcal{I}(0)$ , one should also check directly how well the resulting  $\mathcal{I}(t)$  fits the data.

## More Sophisticated Models

- Let's re-examine the assumptions behind our first models and discuss how to make them more realistic.
- We assumed a fixed population size  $N$  that was isolated from other sources of the hypothetical illness we modeled.
- We assumed that a single number  $p$  represents the probability of an infectious person infecting a susceptible person on each day, for each such pair of people.
- A more realistic model would allow  $p$  to depend on a number of factors.

# Modeling the Infection Probability $p$

- In real life, the infection probability  $p$  depends on the pair of people. However, introducing an independent probability  $p_{mn}$  for each pair of people  $m$  and  $n$  results in way too many parameters.
- Also,  $p$  depends on time; for example, day of week.
- Perhaps most importantly,  $p$  depends on how long the infectious person has had the illness. Typically it peaks a certain amount of time after infection, then decreases to 0.
- To keep the number of parameters manageable, we need to have a model for how  $p$  depends on these factors.



# Compartmental Models

- Many models divide the population into a relatively small number of categories (“compartments”) and keep track of the number of people in each compartment.
- Our first deterministic models had two compartments: “susceptible” and “infectious”. We’ll call the continuous time model (3) the SI model.
- A widely studied model is the SIR model, which introduces a third compartment: “recovered”. People in this category are no longer infectious.
- Other possible compartments can take into account more stages in the progression of the illness, different behavior patterns, different biological characteristics, etc.

# Fitting to Data, Revisited

- In our earlier discussion, we assumed that the number of infectious people at a given time could be measured. But how would we ever know this number?
- The number of infectious people is often inferred from data on new diagnoses of the illness. However:
  - Not all people who get the illness see a doctor.
  - Diagnosis may come well after a person becomes infectious.
  - Data is not always reported (e.g., to CDC) promptly or reliably.
- A common problem in modeling is to relate the quantities of interest to the available data.