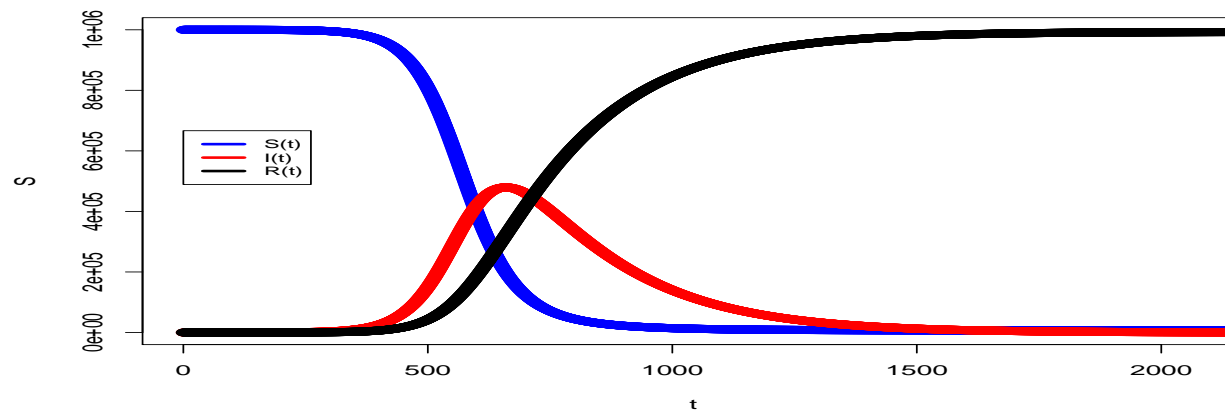


Modeling and Estimation of COVID-19 Under-Reported Counts

Eric Baron¹ and Michael Baron²

¹University of Connecticut, Storrs CT, and ²American University, Washington DC



Supported by the U.S. National Science Foundation

Data Science Conference on COVID-19 [DSCC-19], August 28, 2020

Under-Reported Counts

- ▶ Modeling
 - ▶ Untested people
 - ▶ Unconfirmed transmissions
 - ▶ Under-reported recoveries
 - ▶ Parameters and their dynamics
- ▶ Estimation
 - ▶ Bayesian approach
 - ▶ MCMC
- ▶ Discussion

Epidemic modeling

SIR epidemic model

Population consists of 3 groups: **S**usceptible, **I**nfected, **R**ecovered.

Any individual can move through the states in order, $S \rightarrow I \rightarrow R$



$S_t = \#$ of individuals not infected at time t

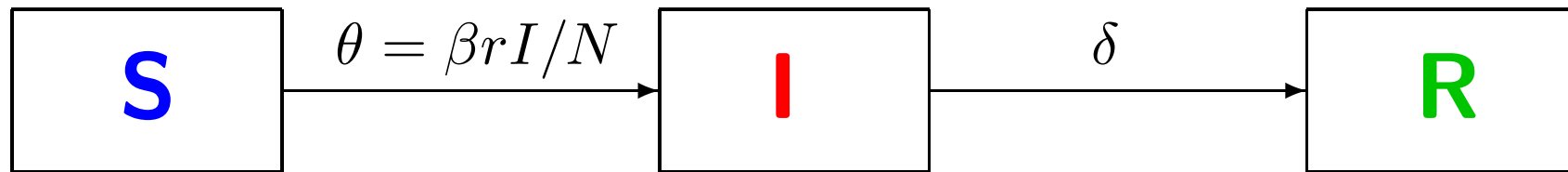
$I_t = \#$ of individuals infected

$R_t = \#$ of individuals infected and then recovered

- ▶ Infection rate θ
- ▶ Recovery rate δ

Kermack and McKendrick 1927

SIR Model



$$\left\{ \begin{array}{l} \frac{dS(t)}{dt} = -\beta r \frac{I(t)}{N} \\ \frac{dI(t)}{dt} = \beta r \frac{I(t)}{N} - \delta I(t) \\ \frac{dR(t)}{dt} = \delta I(t) \end{array} \right. \quad \text{with} \quad \left\{ \begin{array}{l} S(0) = N \\ I(0) = 1 \\ R(0) = 0 \end{array} \right.$$

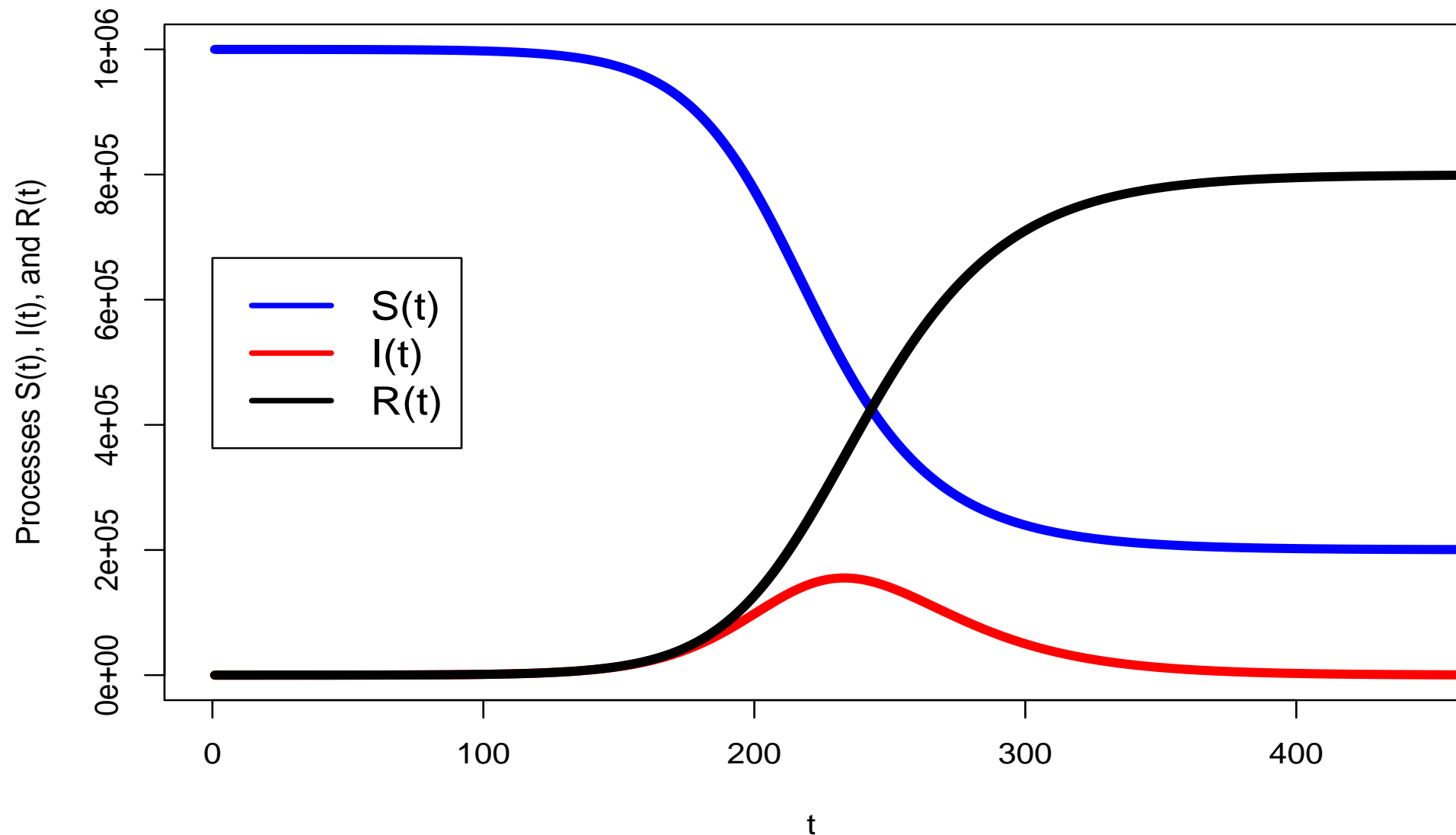
r = number of contacts per unit of time

N = population size

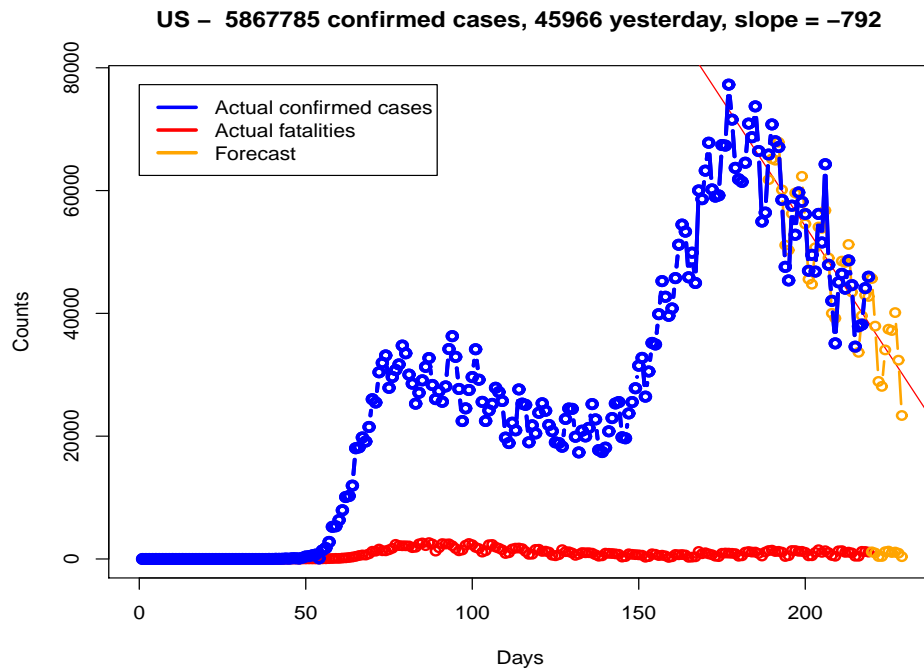
$$\frac{dS}{dt} + \frac{dI}{dt} + \frac{dR}{dt} = 0$$

SIR Model

$\beta = 0.05$, $\delta = 0.05$, $r = 2$, $l_0 = 10$, $N = 1e+06$, $R/N = 80\%$



SIR Model. Does not quite fit the USA data...



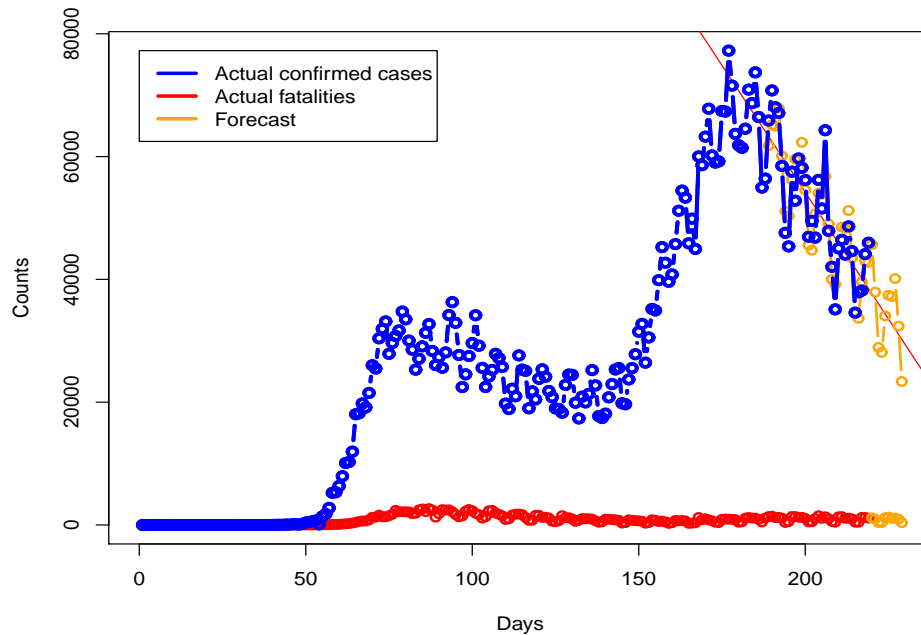
Data from:

Humanitarian Data Exchange (HDX)

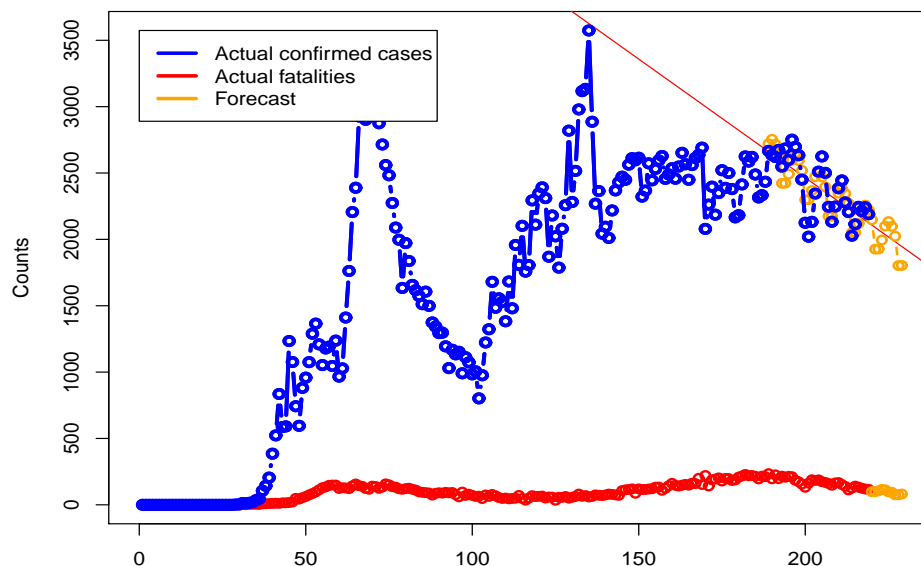
<https://data.humdata.org/dataset/>

SIR Model. Does not quite fit USA, Iran...

US – 5867785 confirmed cases, 45966 yesterday, slope = -792



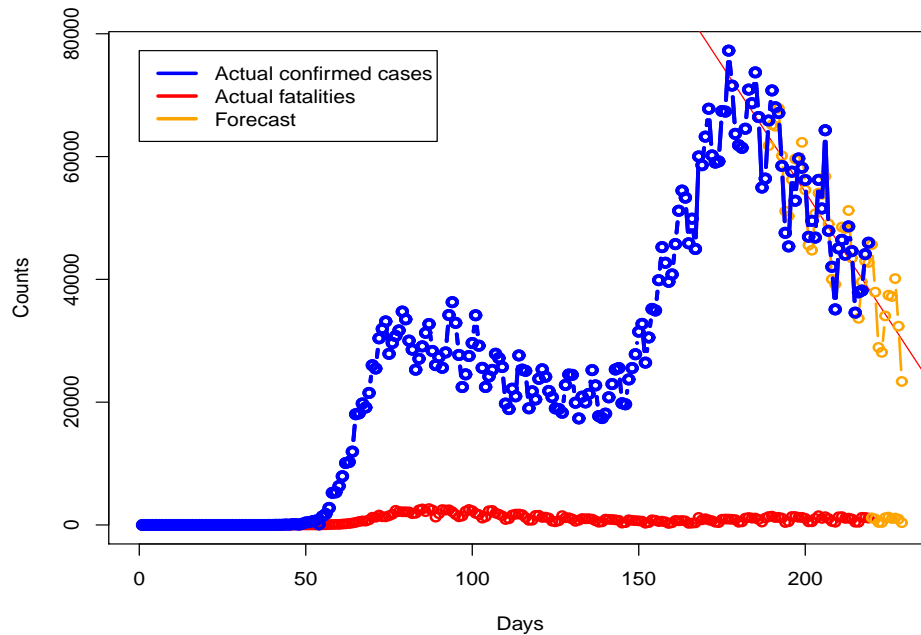
Iran – 367796 confirmed cases, 2190 yesterday, slope = -18



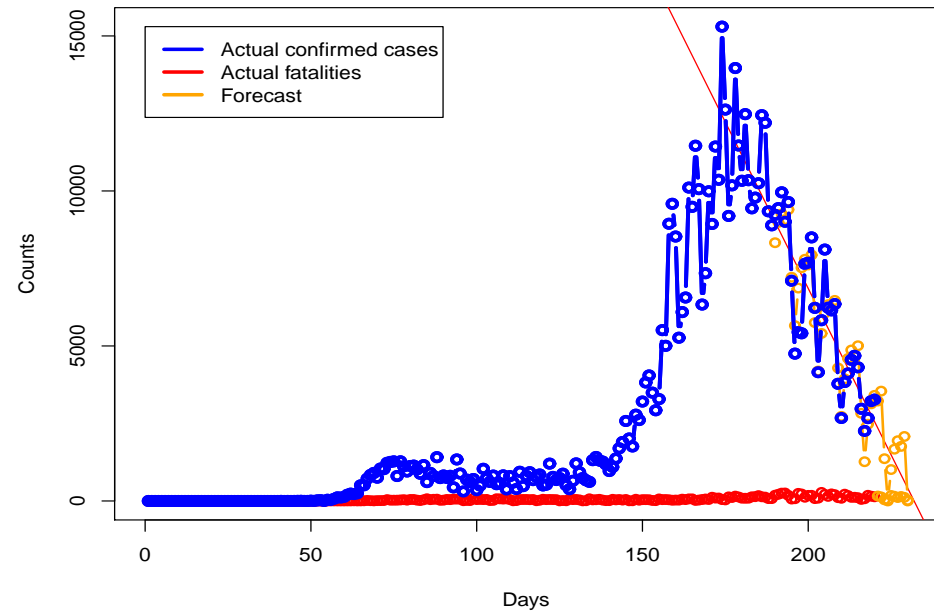
Data from:
Humanitarian Data Exchange (HDX)
<https://data.humdata.org/dataset/>

SIR Model. Does not quite fit USA, Iran, Florida...

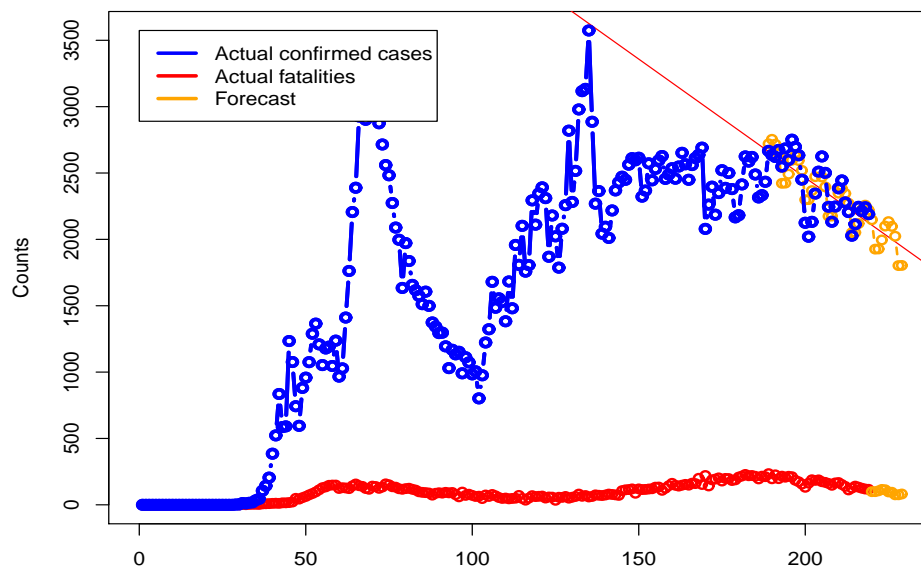
US – 5867785 confirmed cases, 45966 yesterday, slope = -792



Florida – 611983 confirmed cases, 3269 yesterday, slope = -209



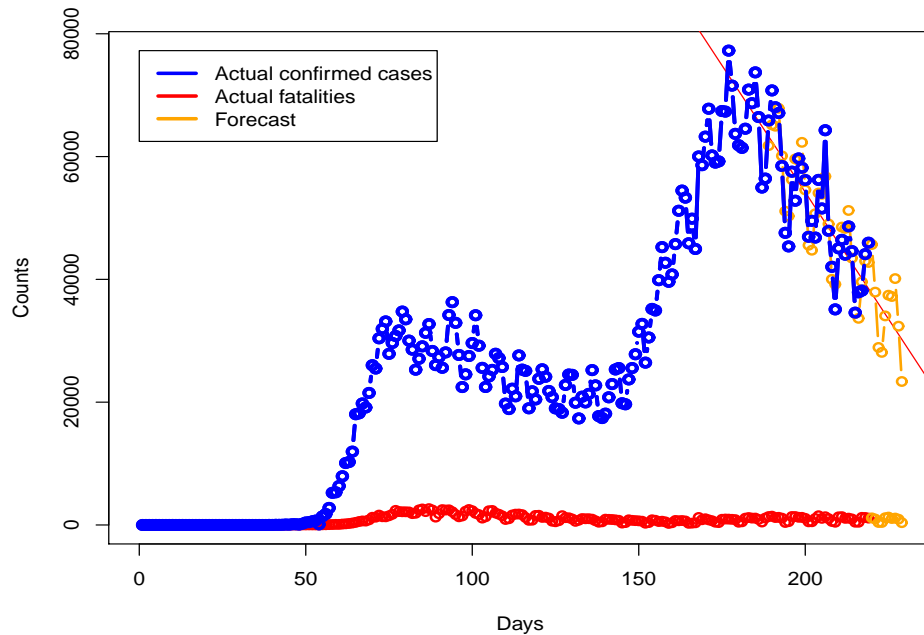
Iran – 367796 confirmed cases, 2190 yesterday, slope = -18



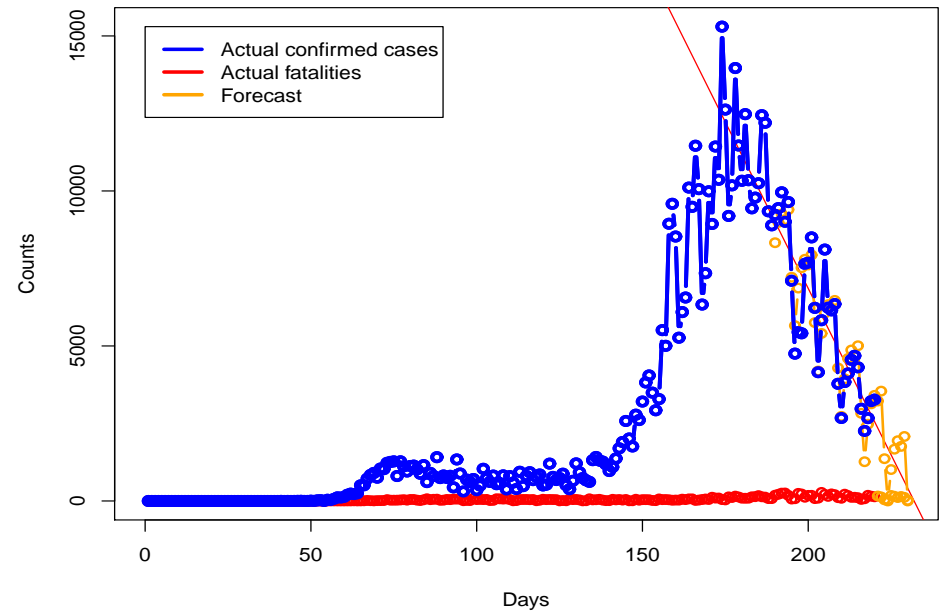
Data from:
Humanitarian Data Exchange (HDX)
<https://data.humdata.org/dataset/>

SIR Model. Does not quite fit USA, Iran, FL, M-D county

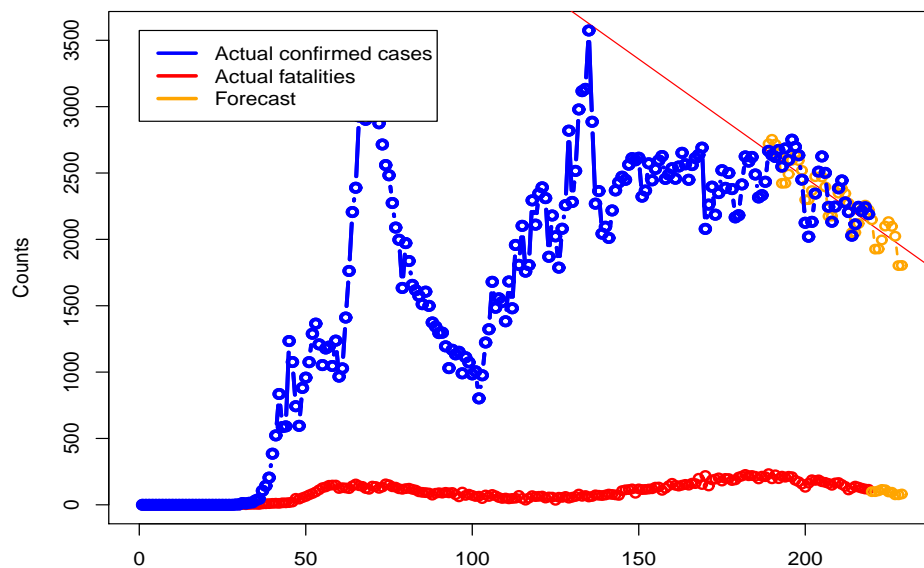
US – 5867785 confirmed cases, 45966 yesterday, slope = -792



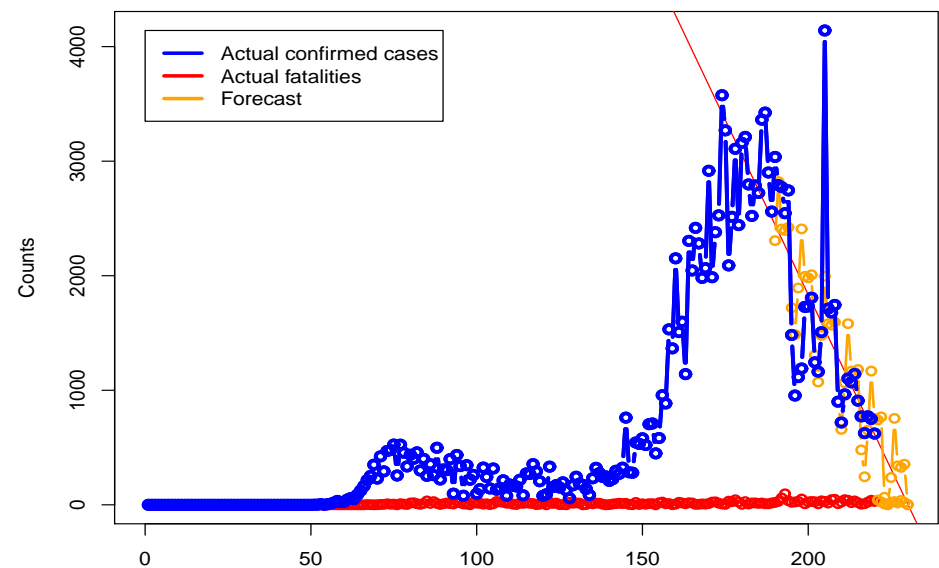
Florida – 611983 confirmed cases, 3269 yesterday, slope = -209



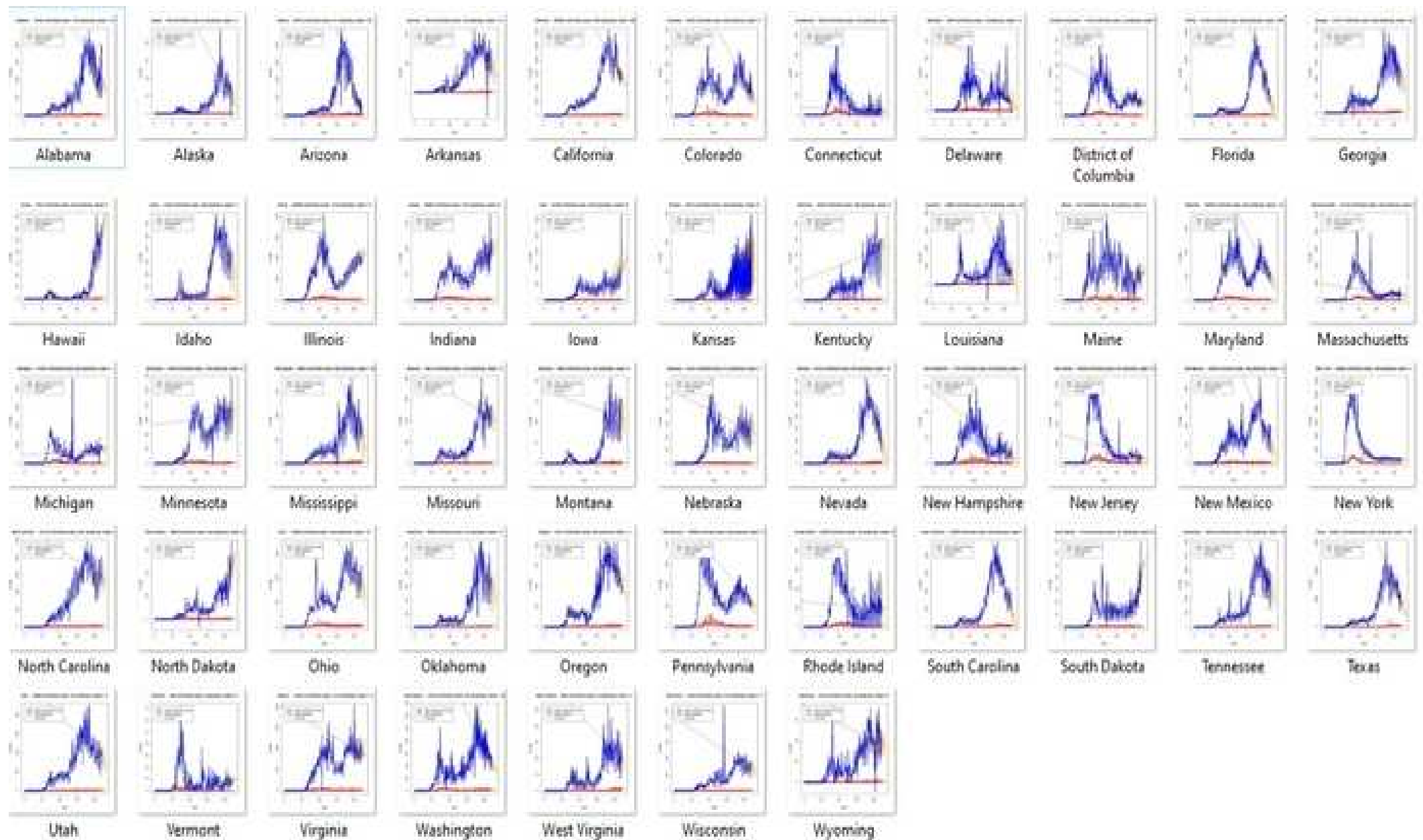
Iran – 367796 confirmed cases, 2190 yesterday, slope = -18



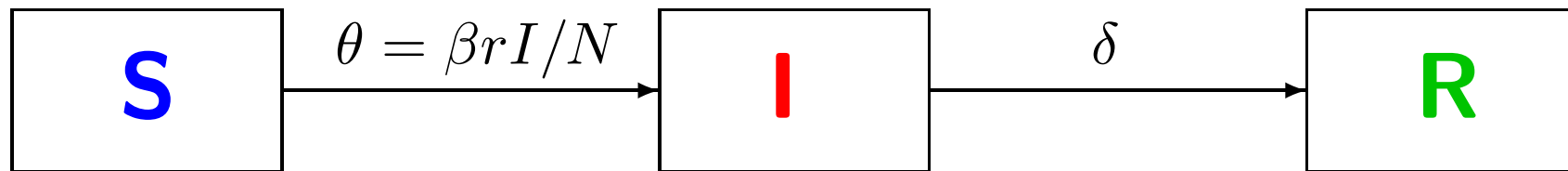
Miami-Dade County – 154756 confirmed cases, 622 yesterday, slope = -59



Why are the states so different?



SIR model



θ = infection rate

β = transmission rate

δ = recovery rate

r = average number of contacts

N = population size

Counts :

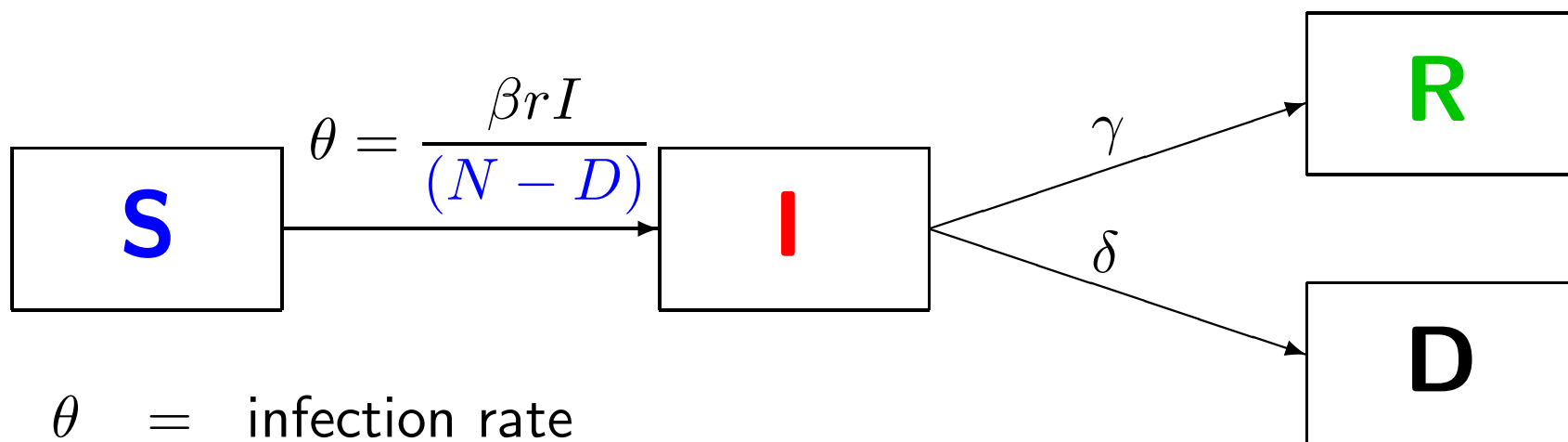
Susceptible

Infected

Recovered

Problem 1: some people do not recover

SIRD model



θ = infection rate

β = transmission rate

γ = recovery rate

δ = mortality rate

r = average number of contacts

N = population size

Counts :

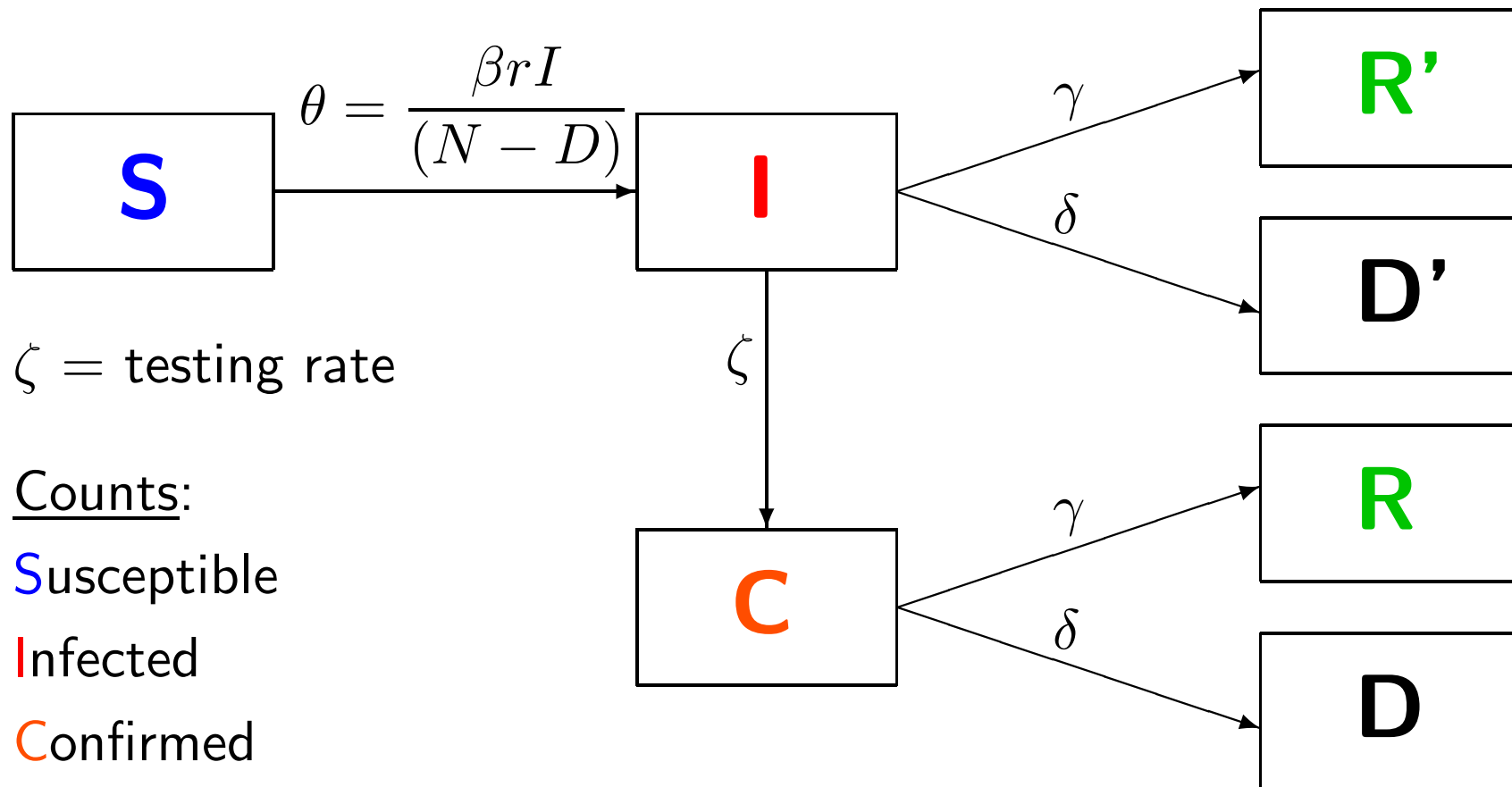
Susceptible Recovered

Infected Died

Problem 2: many people are not tested.

USA: Population 331M; 75M tested; 5.9M confirmed

SICRD model



ζ = testing rate

Counts:

Susceptible

Infected

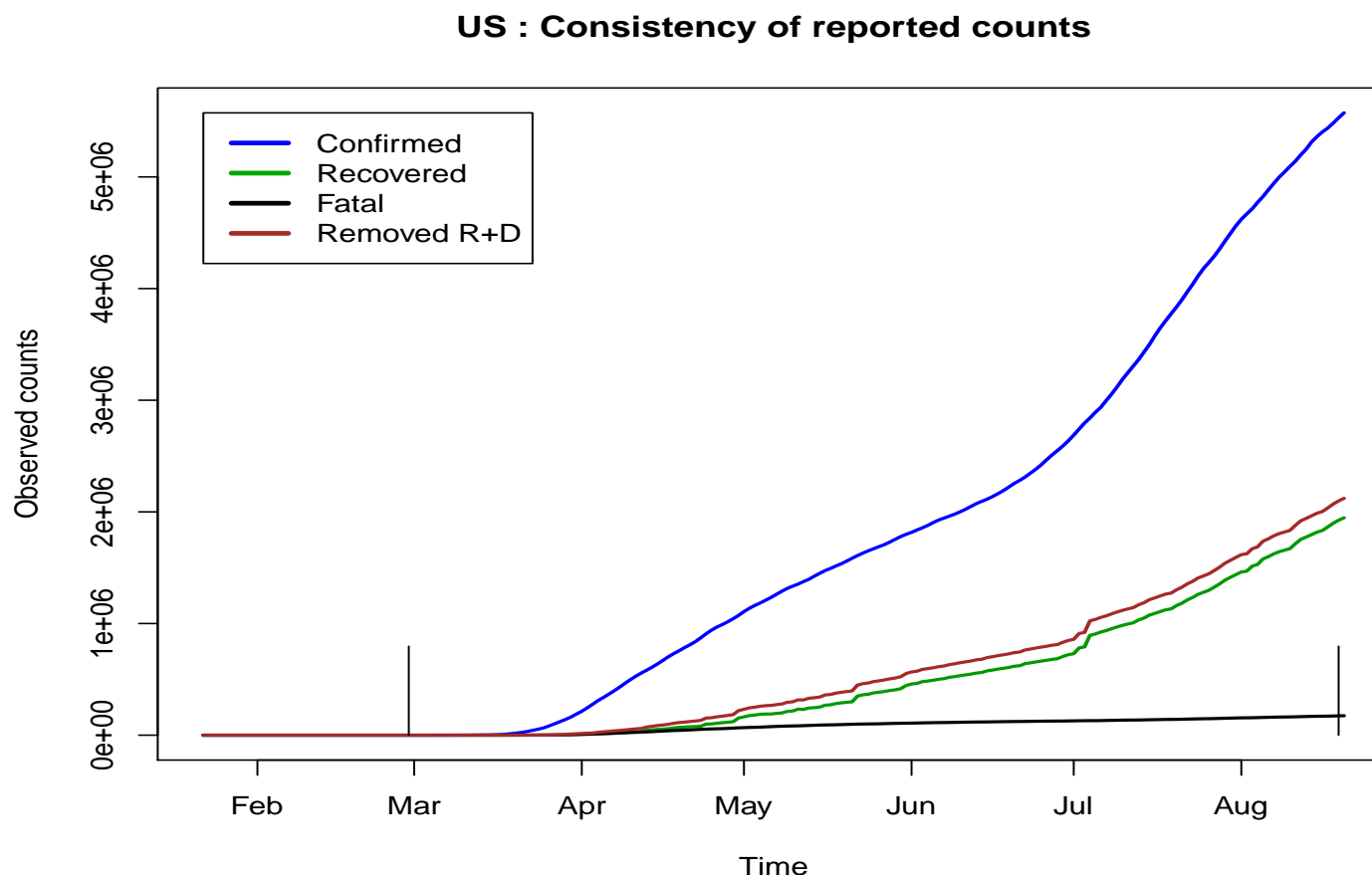
Confirmed

Recovered

Died

Problem 3: under-reported counts

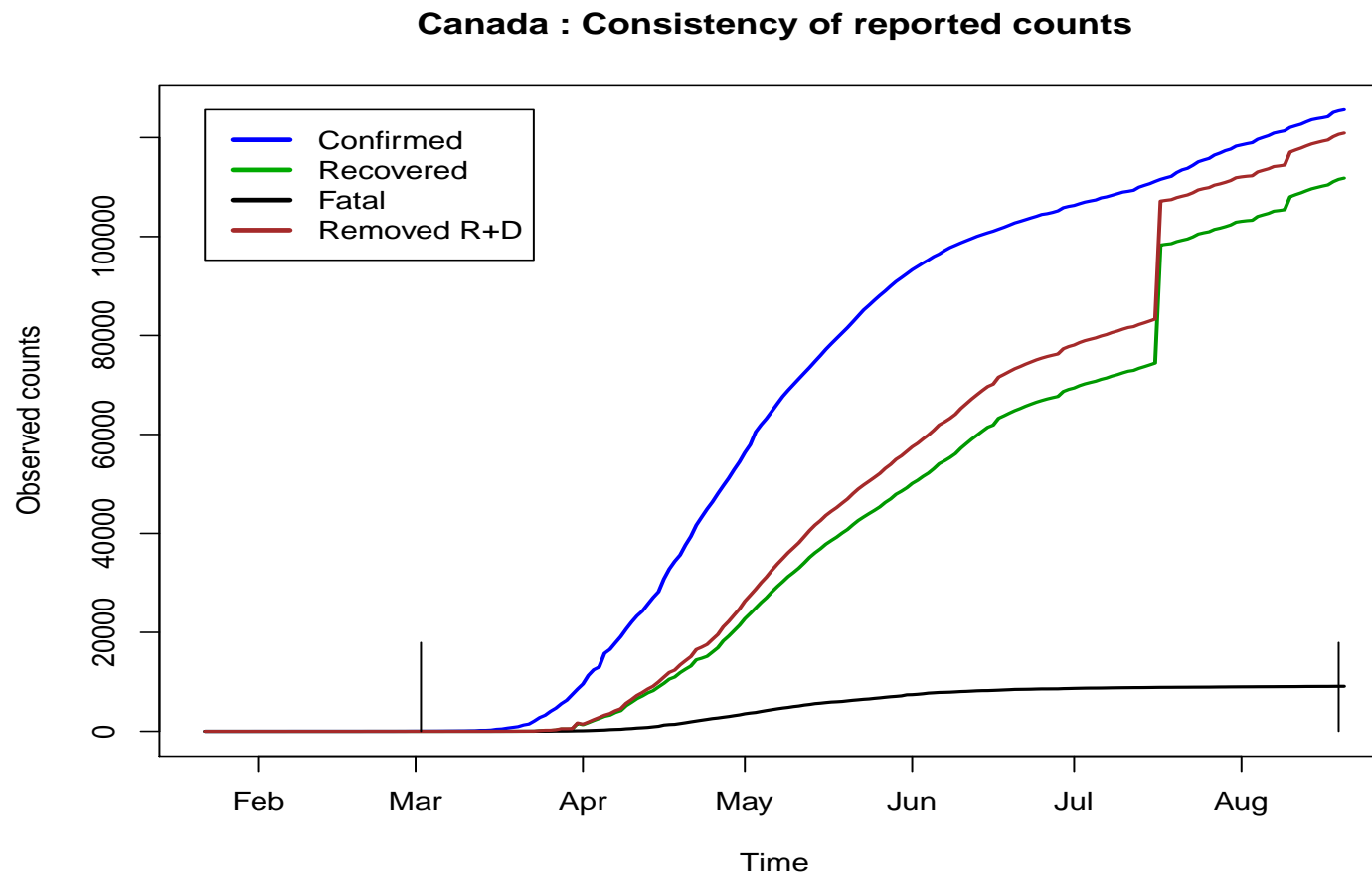
Inconsistency of reported counts: USA



Officially: 2.69M reported confirmed cases in the USA by July 1st
1.95M recovered and 0.17M perished by August 20th

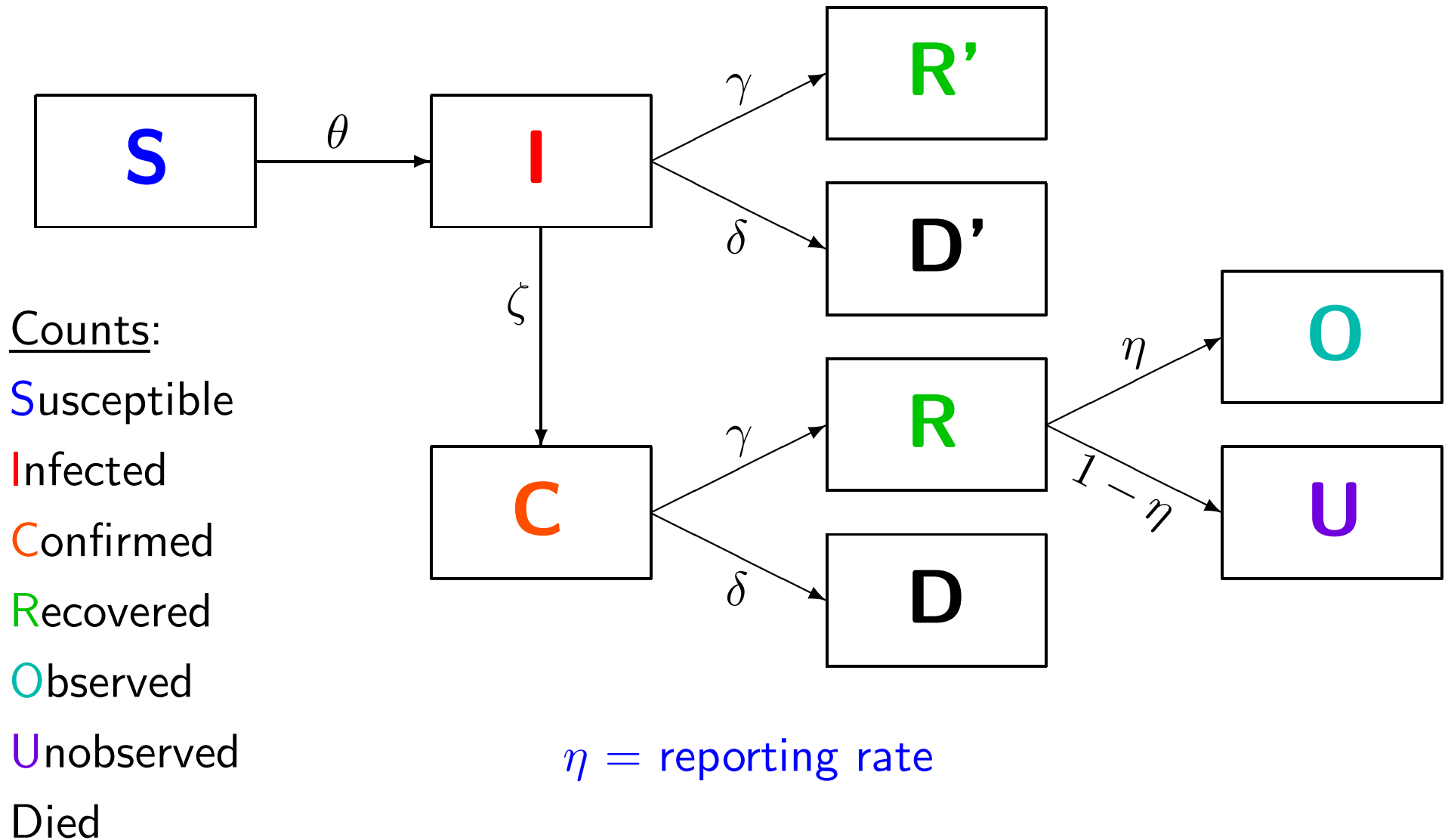
566,298 not reported (actually, many more)

Inconsistency of reported counts: Canada

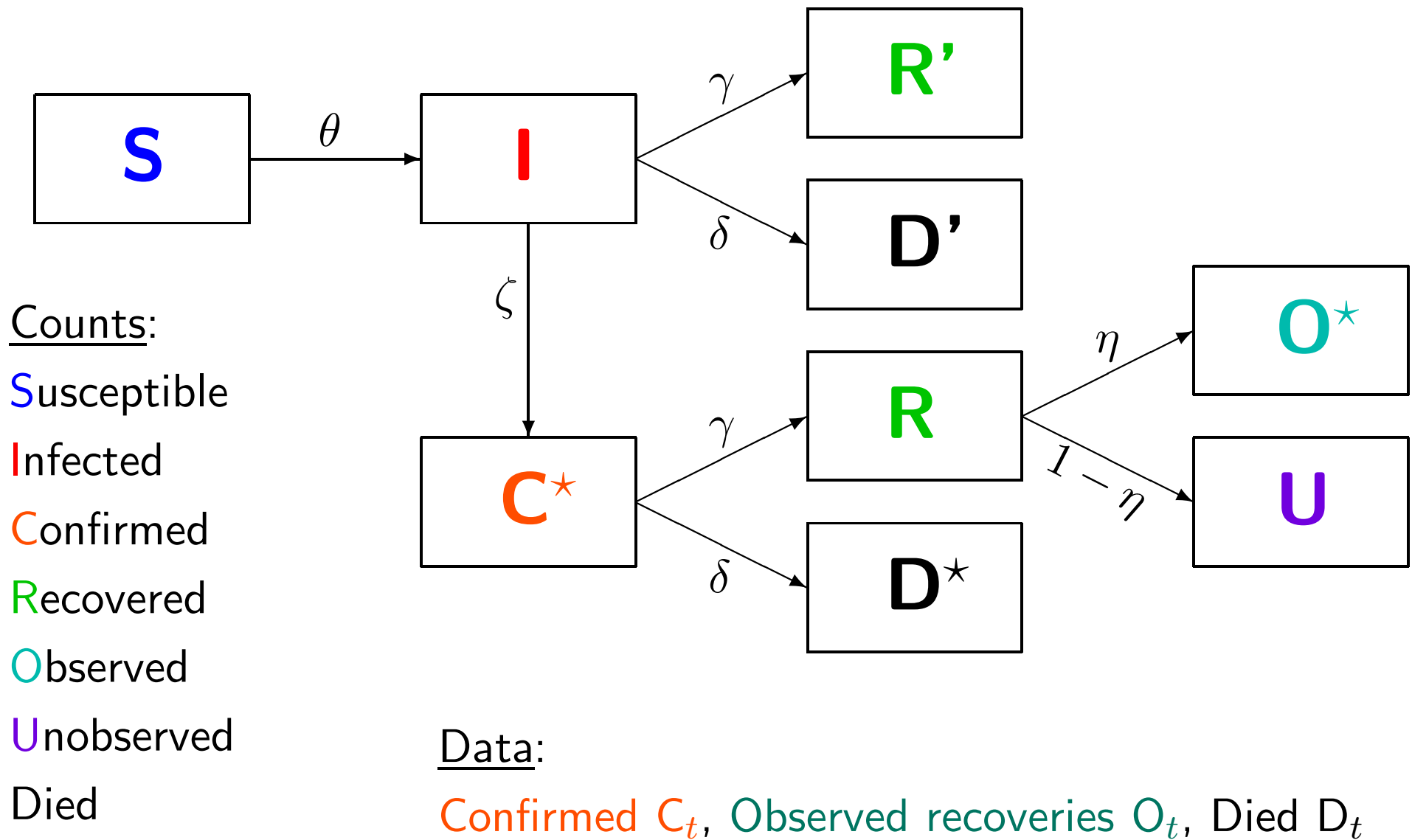


Officially: 23,848 people recovered on July 17
No more than 2,630 people recovered on any other day

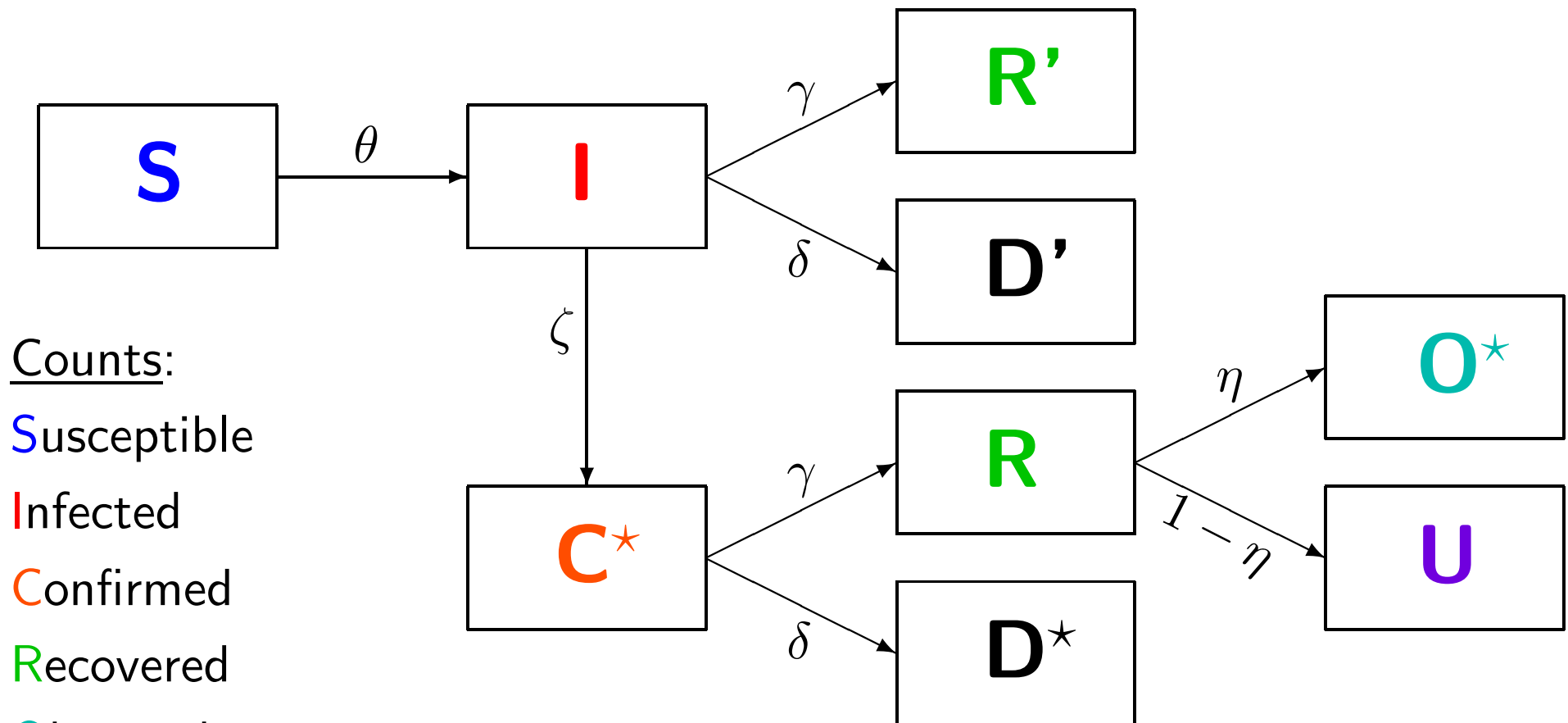
SICROUD model



SICROUD model



SICROUD model



Counts:

Susceptible

Infected

Confirmed

Recovered

Observed

Unobserved

Died

Parameters: $\theta, \gamma, \delta, \zeta, \eta$

Problem 4: SIR model is deterministic.

Bayesian Modeling: Counts and Data

$$|\Delta S_t| \sim \text{Binomial}(S_t, \theta_t)$$

$$\Delta R_t \sim \text{Binomial}(C_t, \gamma_t)$$

$$\Delta D_t \sim \text{Binomial}(C_t, \delta_t)$$

$$\Delta R'_t \sim \text{Binomial}(I_t, \gamma_t)$$

$$\Delta D'_t \sim \text{Binomial}(I_t, \delta_t)$$

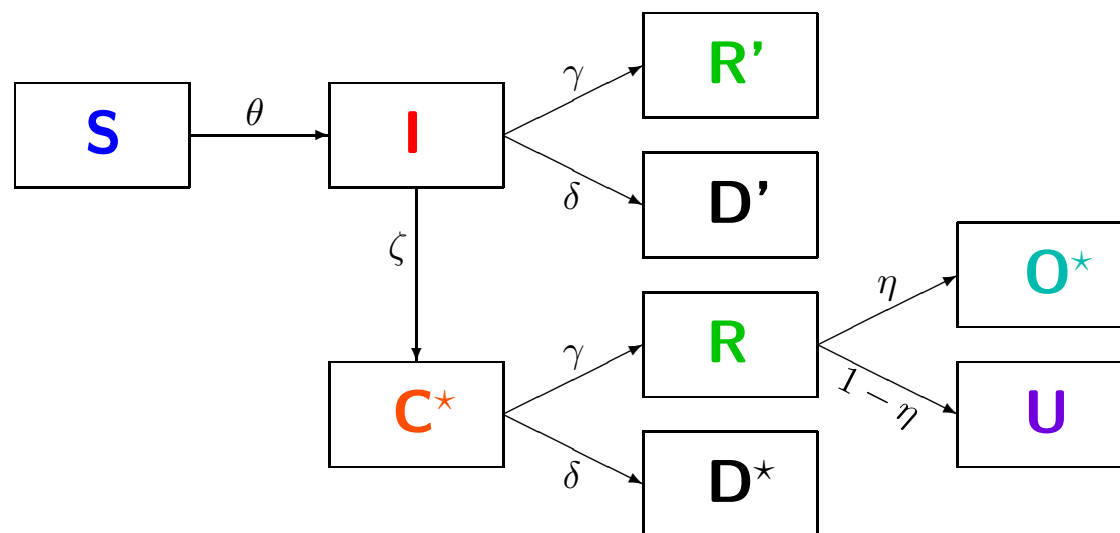
$$\Delta O_t \sim \text{Binomial}(R_t, \eta_t)$$

$$\Delta U_t = R_t - \Delta O_t$$

$$\Delta C_t \sim \text{Binomial}(I_t, \zeta_t) - \Delta R_t - \Delta D_t$$

$$\Delta I_t = |\Delta S_t| - \Delta C_t - \Delta R'_t - \Delta D'_t$$

$$I_t \sim \text{Poisson}(\rho_t)$$

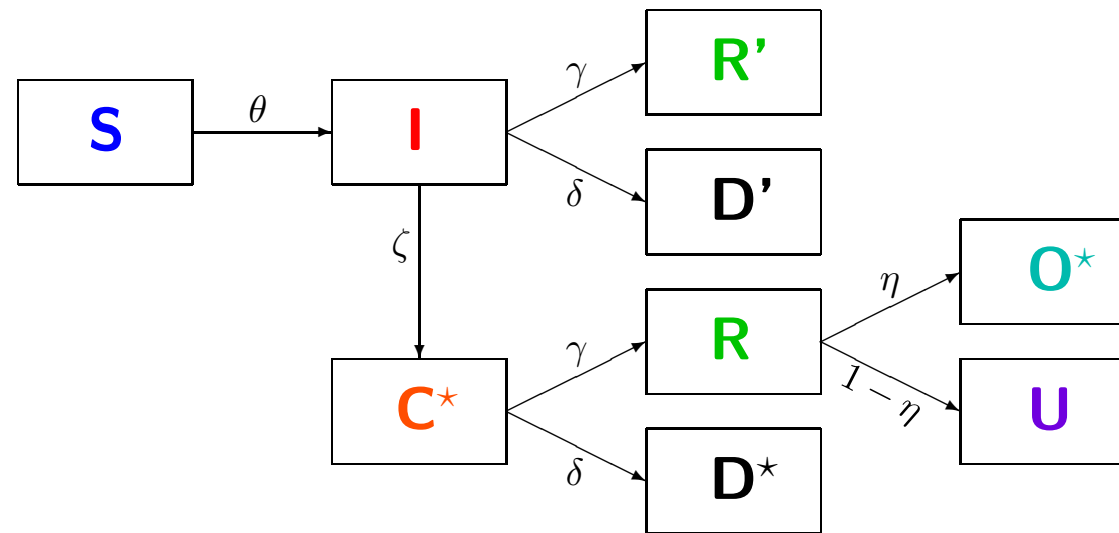


Problem 5:

Parameters are not constant.

$C_t, O_t, D_t, \Delta C_t, \Delta O_t, \Delta D_t$ are observed

Bayesian Modeling: Parameters

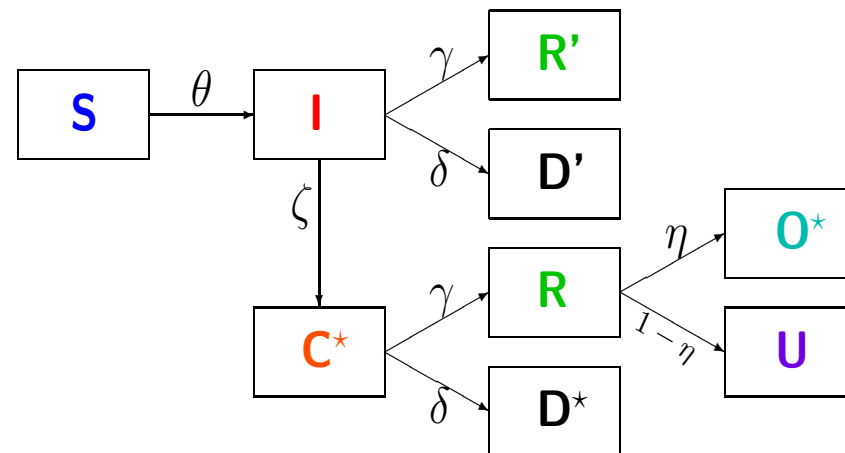


Daily rates:

Infection rate	θ_t	\sim	$\text{Beta}(\alpha^{(\theta)}, \beta^{(\theta)})$	} Prior distributions
Recovery rate	γ_t	\sim	$\text{Beta}(\alpha^{(\gamma)}, \beta^{(\gamma)})$	
Mortality rate	δ_t	\sim	$\text{Beta}(\alpha^{(\delta)}, \beta^{(\delta)})$	
Testing rate	ζ_t	\sim	$\text{Beta}(\alpha^{(\zeta)}, \beta^{(\zeta)})$	
Reporting rate	η_t	\sim	$\text{Beta}(\alpha^{(\eta)}, \beta^{(\eta)})$	
Currently infected	ρ_t	\sim	$\text{Gamma}(\alpha^{(\rho)}, \lambda^{(\rho)})$	

Posterior Distributions: Update of Hyper-Parameters

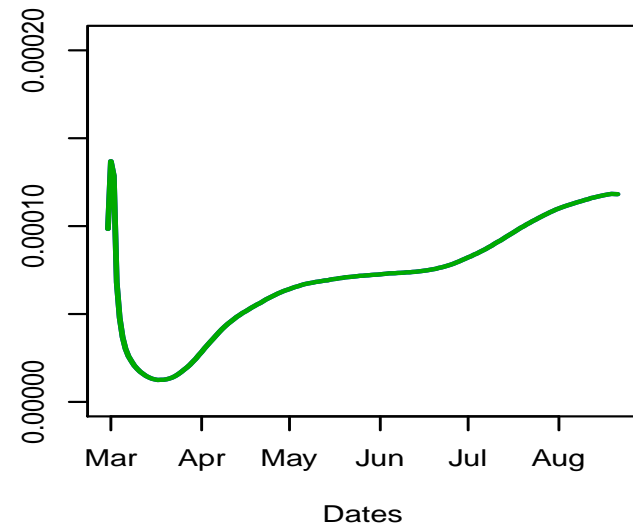
$$\left\{ \begin{array}{l} \theta | I \sim \text{Beta}(\alpha^{(\theta)} + I, \beta^{(\theta)} + S - I) \\ \gamma | R \sim \text{Beta}(\alpha^{(\gamma)} + R, \beta^{(\gamma)} + C - R) \\ \delta | D \sim \text{Beta}(\alpha^{(\delta)} + D, \beta^{(\delta)} + C - D) \\ \zeta | C \sim \text{Beta}(\alpha^{(\zeta)} + C, \beta^{(\zeta)} + I - C) \\ \eta | O \sim \text{Beta}(\alpha^{(\eta)} + O, \beta^{(\eta)} + R - O) \\ \rho | I \sim \text{Gamma}(\alpha^{(\rho)} + I, \beta^{(\rho)} + 1) \\ I | C \sim \text{Poisson}(\rho), \text{shifted by } C \end{array} \right.$$



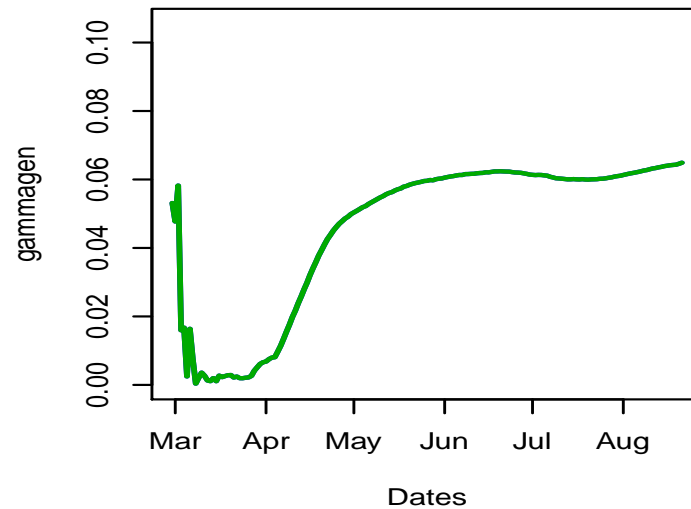
- Dynamics of the COVID-19 pandemic:
- Observed counts induce daily updates of hyper-parameters
- The Beta-Binomial model is ready for forecasting

Parameter Estimates: USA

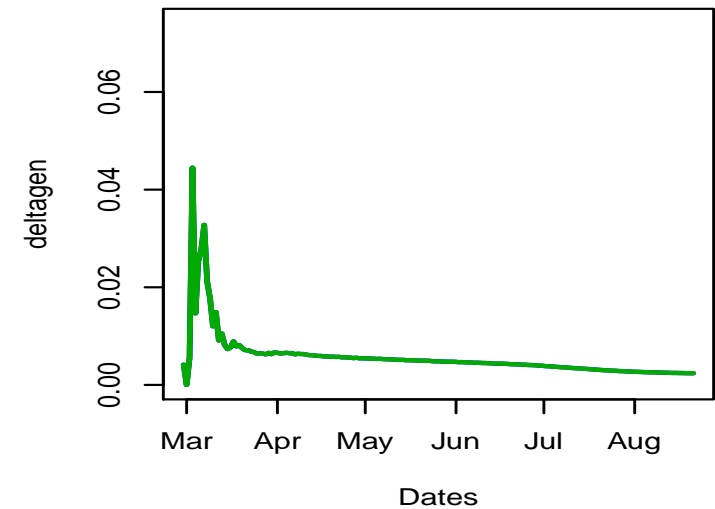
US - ttheta



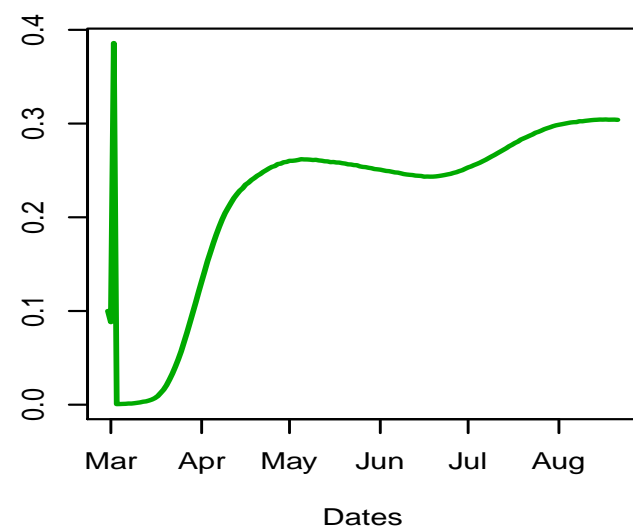
US - gamma



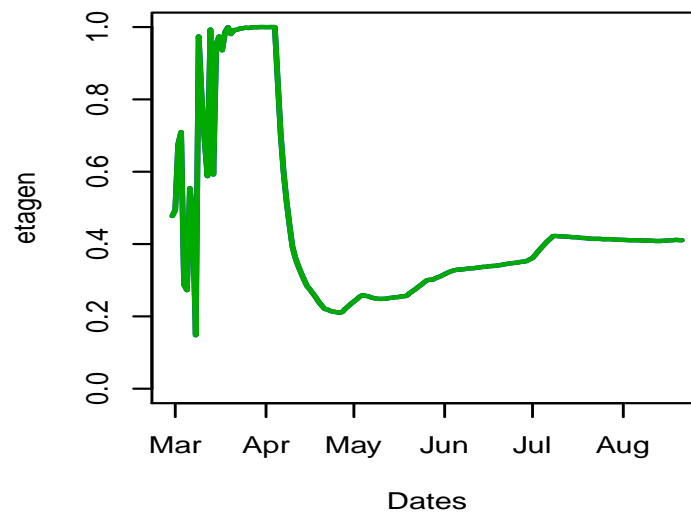
US - delta



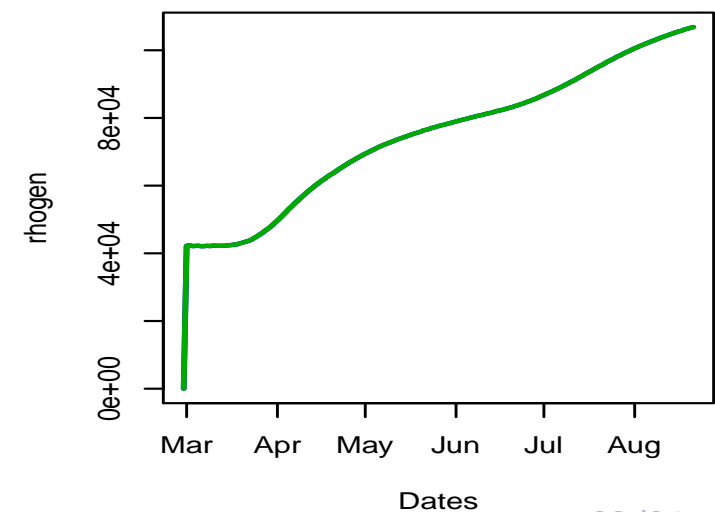
US - zeta



US - eta



US - rho



Latest Estimates: USA (08/28/2020)

Infection rate	θ	=	0.00012
Transmission rate	β	=	0.056
Recovery rate	γ	=	0.066
Mortality rate	δ	=	0.0.0023
Testing rate	ζ	=	0.369
Currently infected, not tested	ρ	=	88,944
Basic reproduction number	R_0	=	0.831
Mean disease duration	$\frac{1}{\gamma+\delta}$	=	14.7 days
Total infected	$\sum \Delta S_t $	=	6,964,213
Confirmed	C	=	5,801,712
Total recovered	$R + R'$	=	6,014,168
Reported recoveries	O	=	2,065,066
Total casualties	$D + D'$	=	258,993
Reported casualties	D	=	179,066

Discussion

- ▶ This approach fills *some* gaps in official reports
- ▶ It brings consistency among observed counts
- ▶ But it may not recover *everything* that is hidden
- ▶ It allows parameters to change, reflecting the epidemic dynamics
- ▶ It provides natural forecasting
- ▶ It can be extended accounting for new types of data

Any questions?

Thank you!

