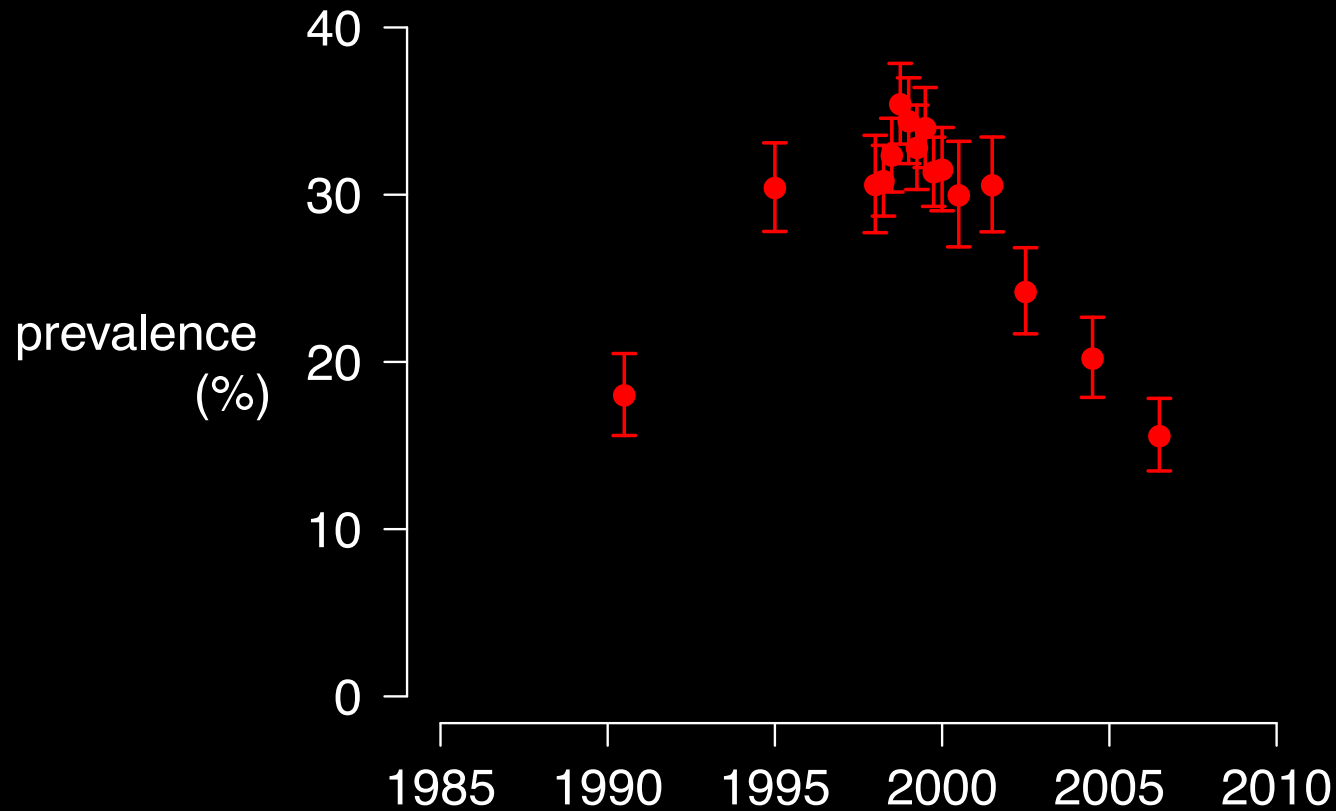# Models and Data

## Introduction to Model Fitting



Steven Bellan, PhD, MPH
Dept. of Epidemiology & Biostatistics, College of Public Health
University of Georgia
DAIDD, White Oak Conservation
Thursday December 8, 2016

1

# Outline

1. Recap: Classical and Mechanistic Epidemiology

2. Why fit models to data?

3. Review of Linear Regression

4. Maximum Likelihood and Fitting Simple Models

5. Fitting Dynamic Models to Data
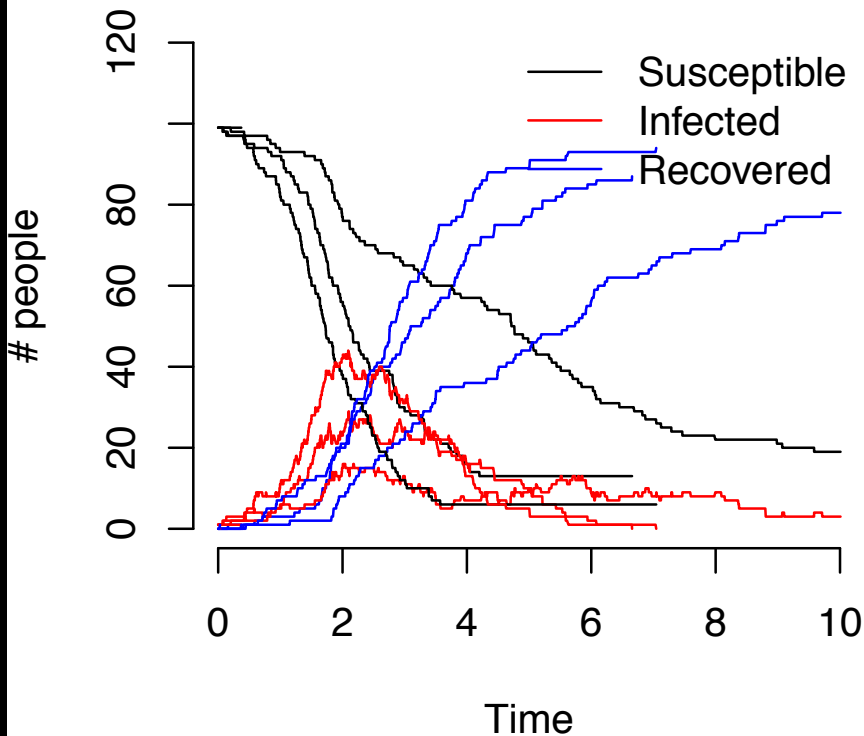
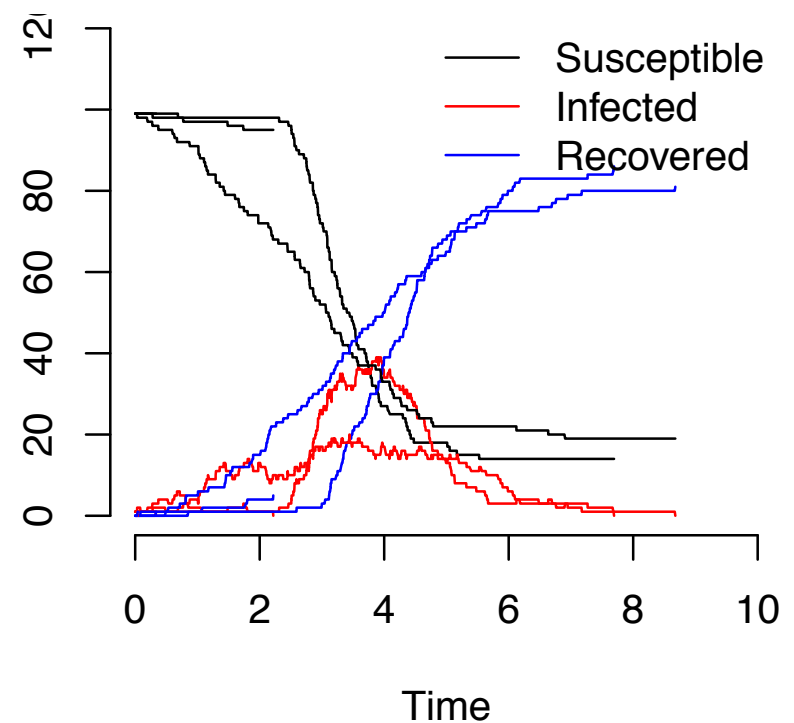6. Summary

# What happened?

## Harare ANC HIV Data

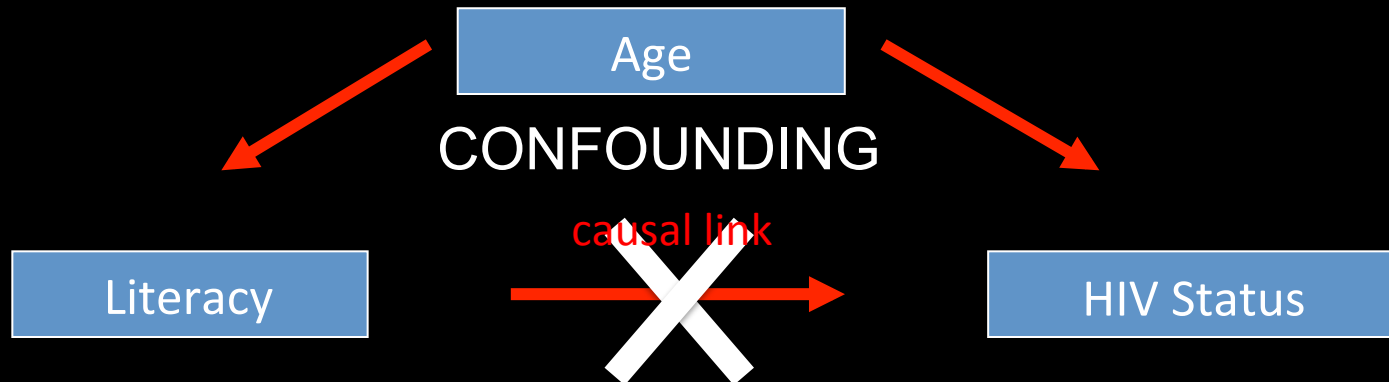# Are these different?

## Measles Outbreaks

# Classical Epidemiology

| Individual | Literate | HIV infected |
|---|---|---|
| 1 | 0 | 0 |
| 2 | 0 | 0 |
| 3 | 0 | 0 |
| 4 | 0 | 1 |
| 5 | 1 | 1 |
| 6 | 1 | 0 |
| 7 | 1 | 1 |
| 8 | 1 | 1 |

- Does literacy cause HIV?

- Find correlations that imply causality by accounting for

  1. random error: do we have enough data?

  2. bias: are design & analysis valid?

Age

CONFOUNDING

causal link

Literacy  ✕  HIV Status
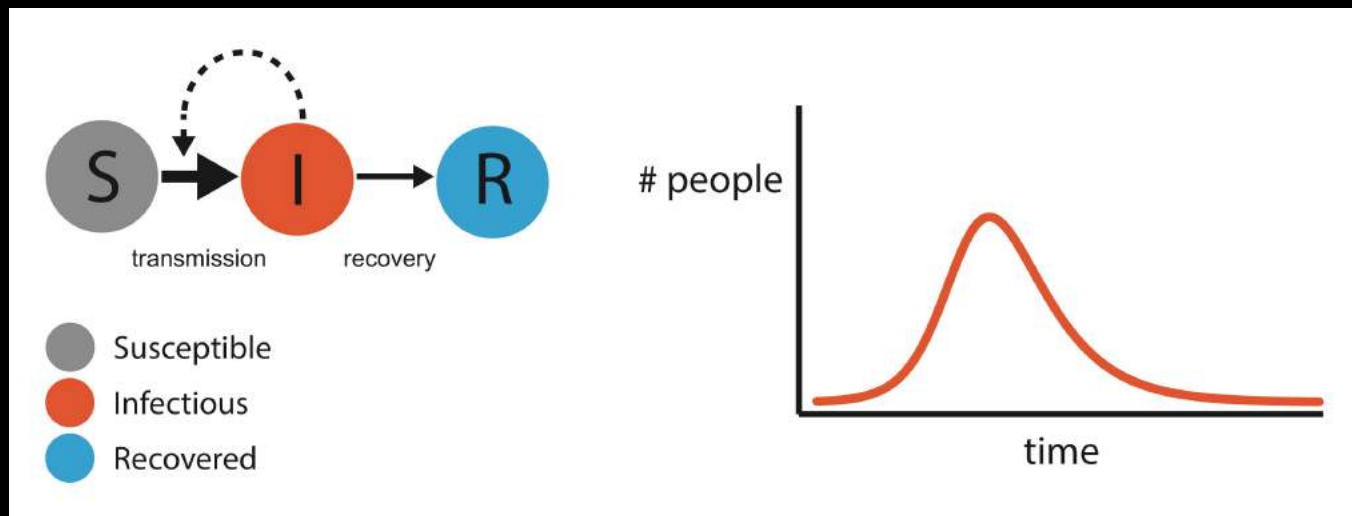
# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

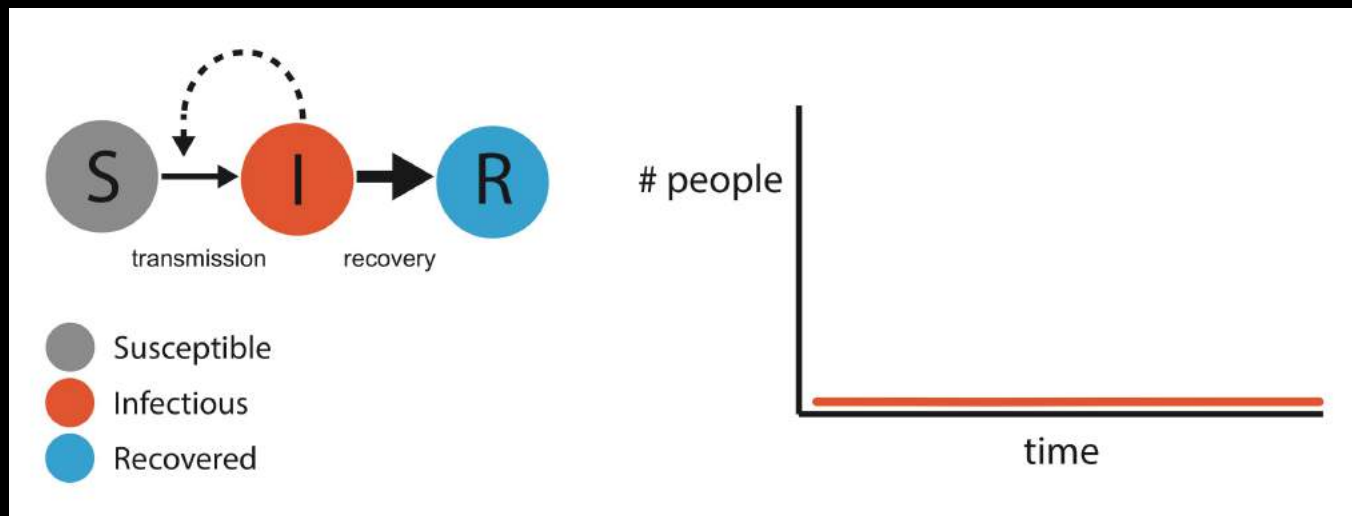- "What if" scenarios not amenable to experimentation

What if each person exposed 50% more people?
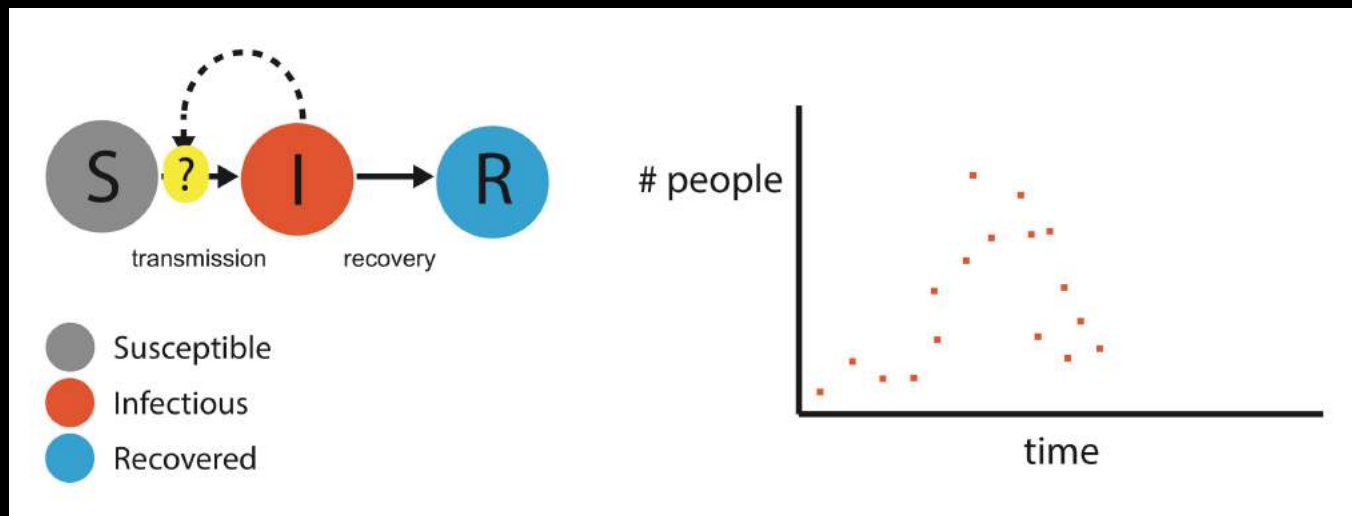
# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

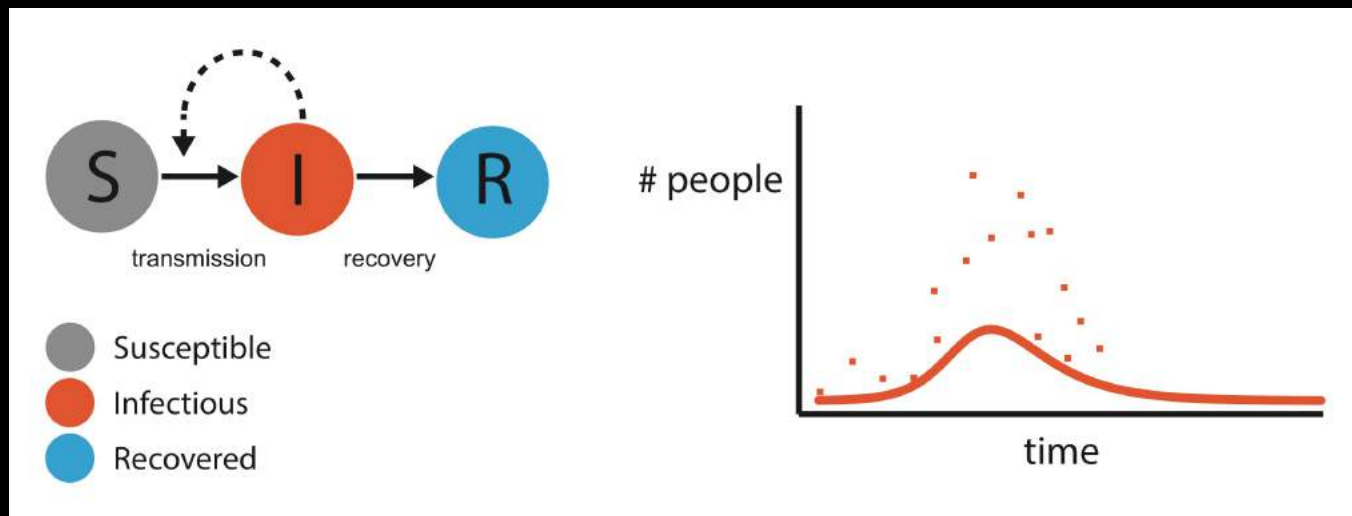What if we treated people and doubled the rate of recovery?

# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

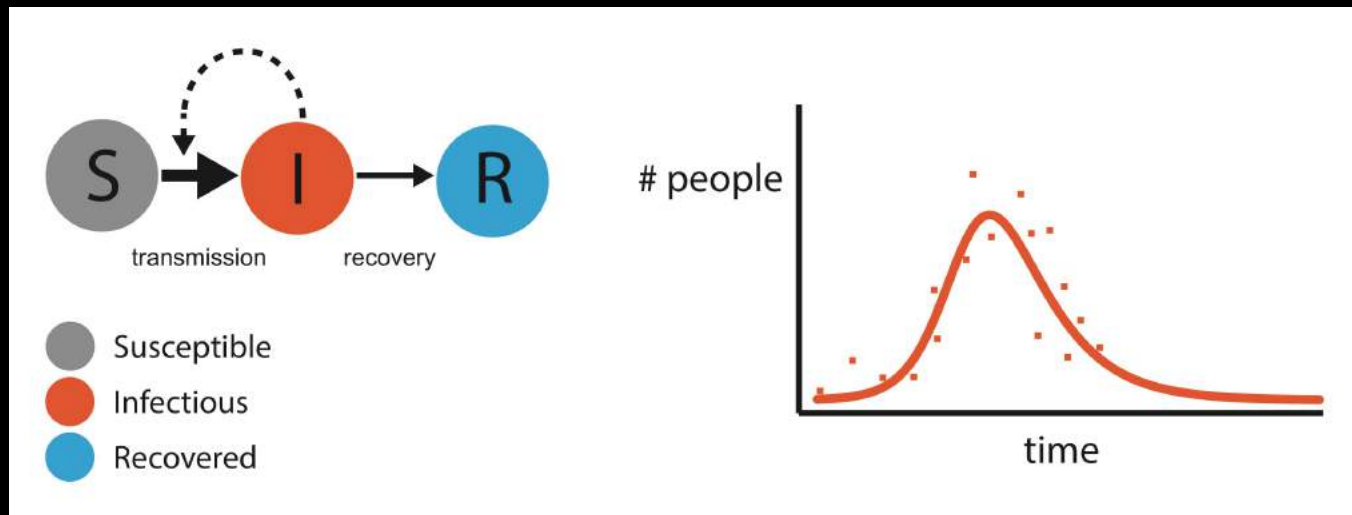- Estimating parameters by fitting available data

# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

- Estimating parameters by fitting available data

# Mechanistic Epidemiology
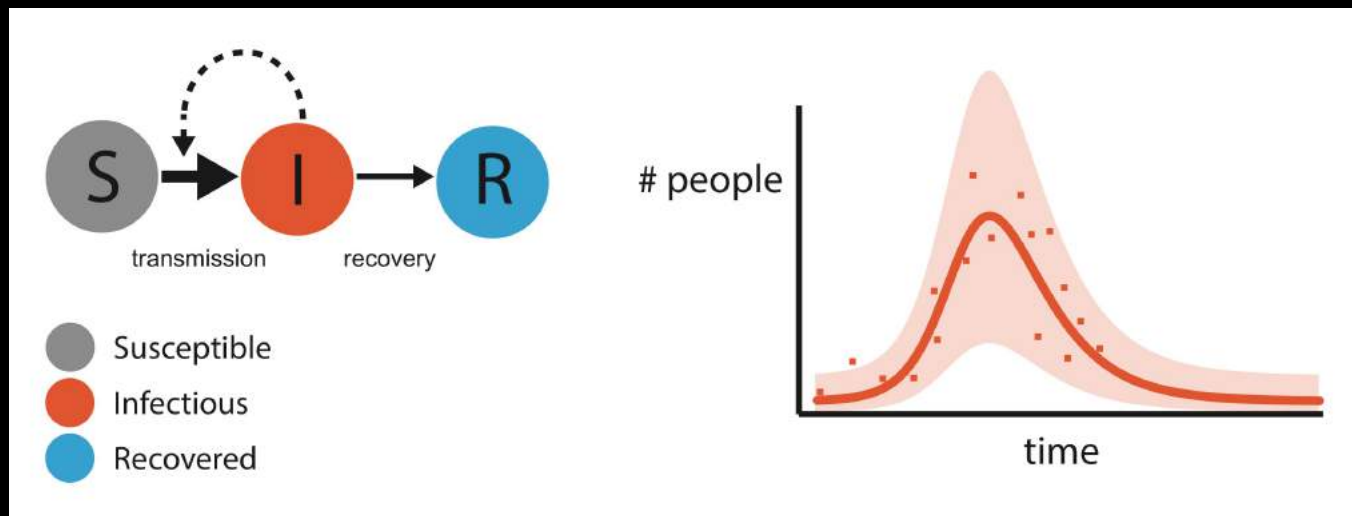
- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

- Estimating parameters by fitting available data

# Mechanistic Epidemiology
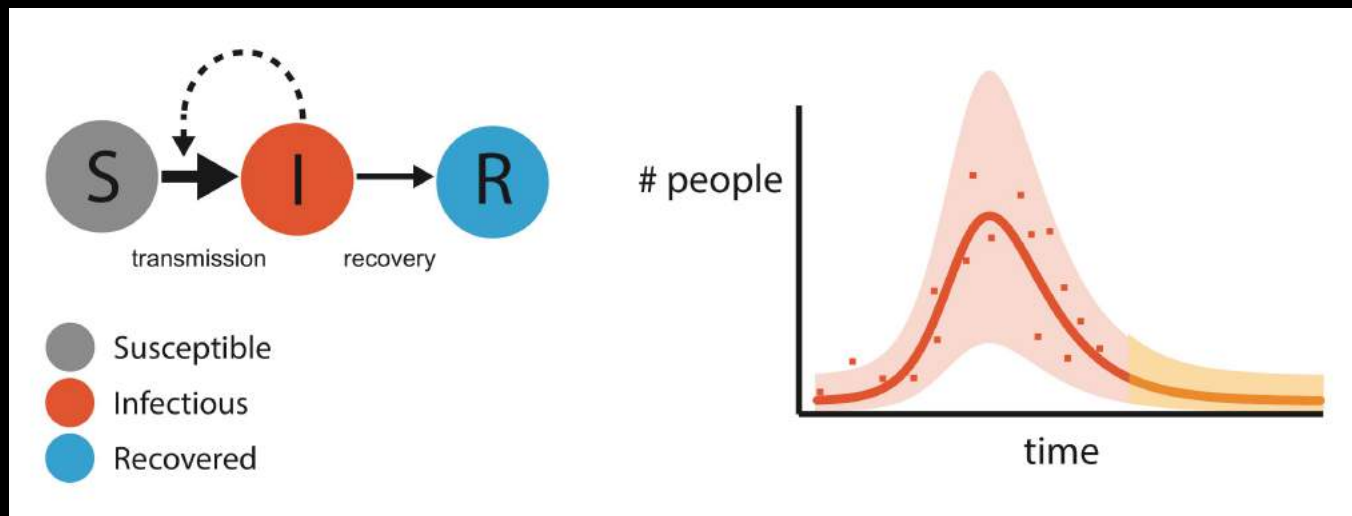
- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

- Estimating parameters by fitting available data

Estimate transmission rate or other model parameters
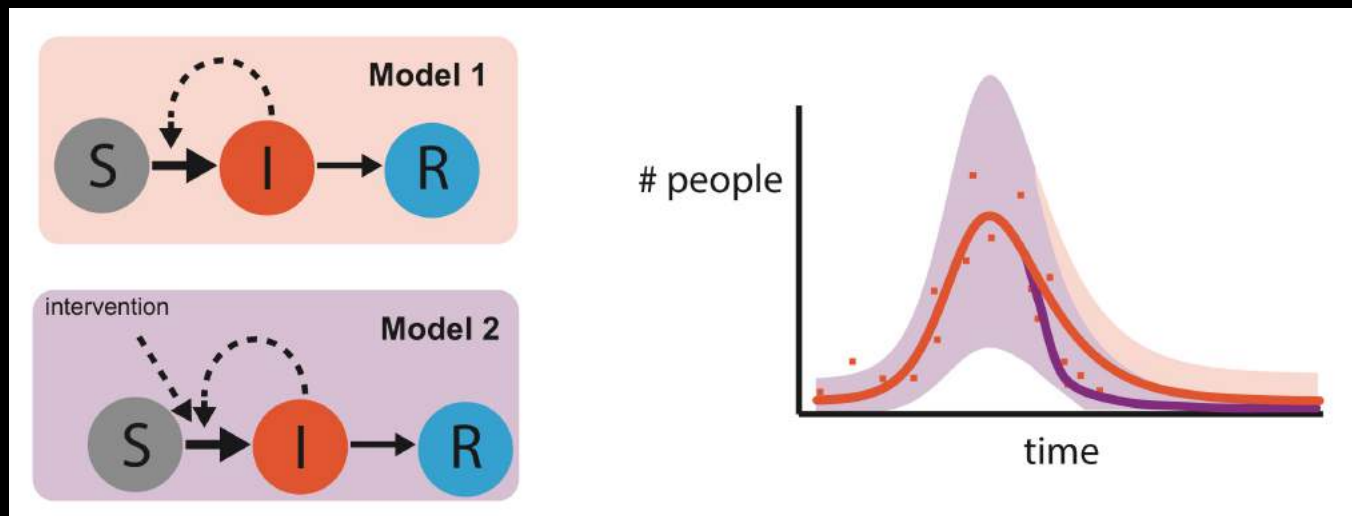(with confidence intervals)

# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

- Estimating parameters by fitting available data
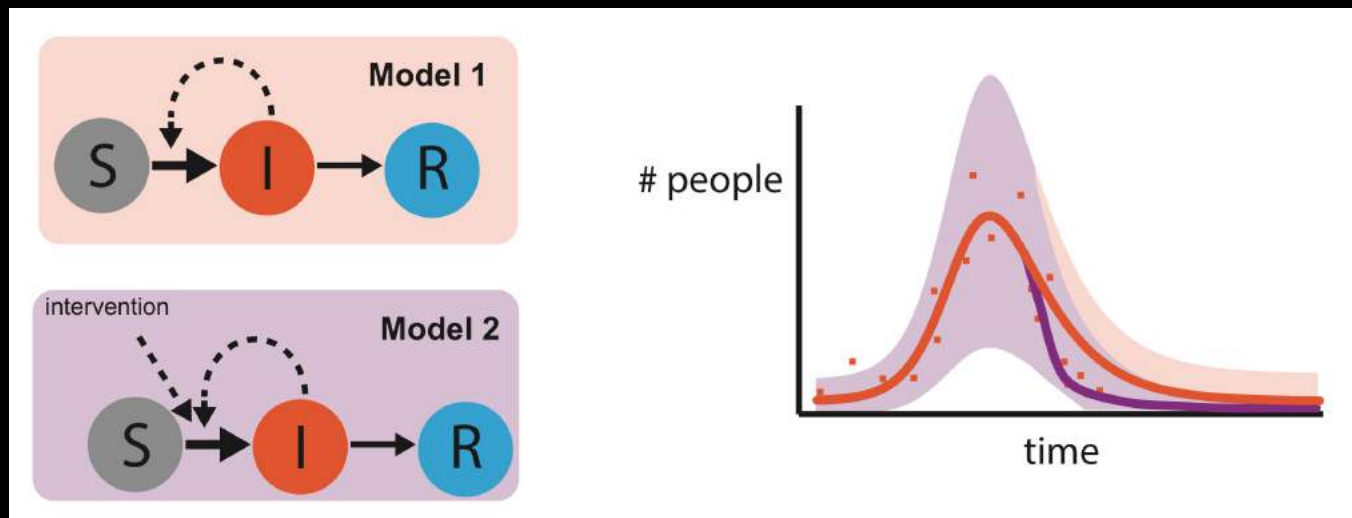
- Prediction



13

# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

- Estimating parameters by fitting available data

- Prediction

- Model selection (choosing between alternative hypotheses)

# Mechanistic Epidemiology

- Scale up from individual processes to population patterns

- "What if" scenarios not amenable to experimentation

- Estimating parameters by fitting available data

- Prediction

- Model selection

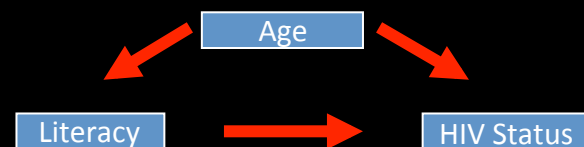data focus emerged in last 10 years

# Outline

1. Recap: Classical and Mechanistic Epidemiology

2. Why fit models to data?

3. Review of Linear Regression

4. Maximum Likelihood and Fitting Simple Models

5. Fitting Dynamic Models to Data

6. Summary

# Why fit models to data?

- **Estimate** quantities/parameters of interest

- **Inference**: Test hypotheses

- Model assessment:

    Assess **plausibility** or **model comparison**

- End goal: **explain** observed patterns or **predict**

# Statistical Models
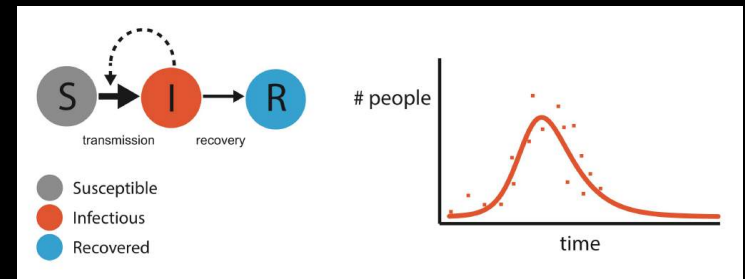
- A **familiar** starting point



- **Analogous** to fitting dynamical models

- **Abstraction** of real relationships

- **Explaining variation** in data through **correlational** relationships (hopefully causal)

# Dynamic Models and Time Series Data

- Dynamic models evolve through time

- and simulate time series



- Informally compare observed time series & simulated time series

- Fitting models to data formally compares them
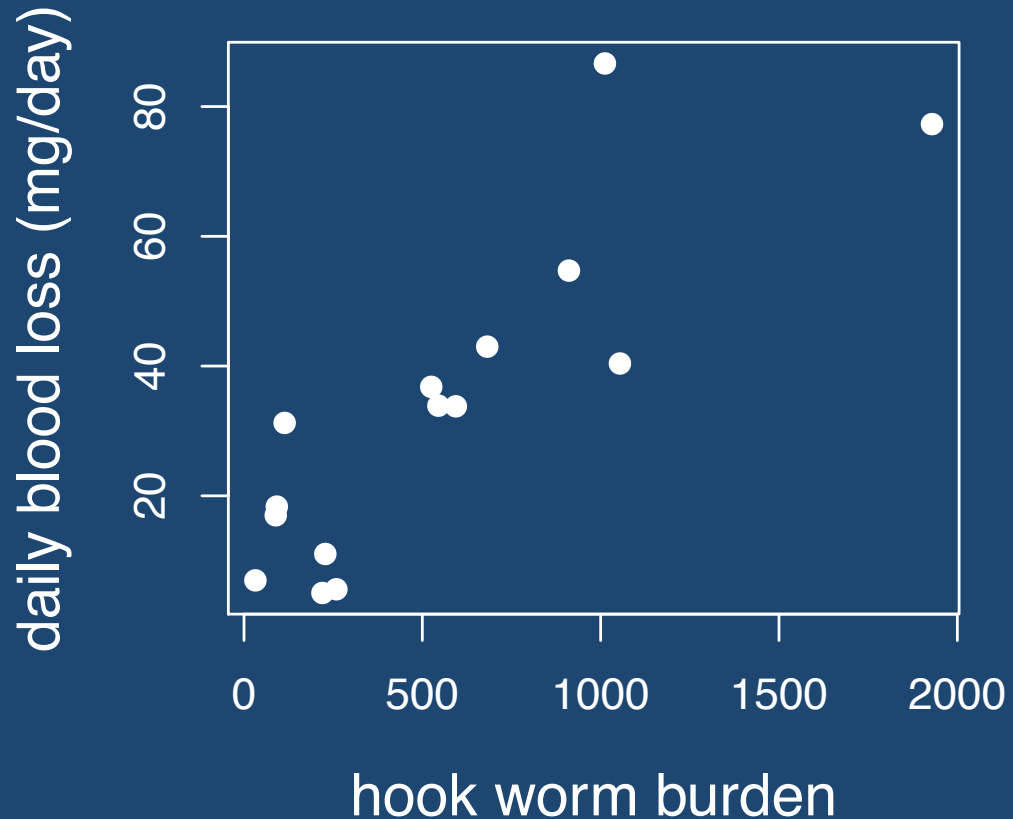
# Outline

1.  Recap: Classical and Mechanistic Epidemiology

2.  Why fit models to data?

3.  Review of Linear Regression

4.  Maximum Likelihood and Fitting Simple Models

5.  Fitting Dynamic Models to Data

6.  Summary

# Linear Regression

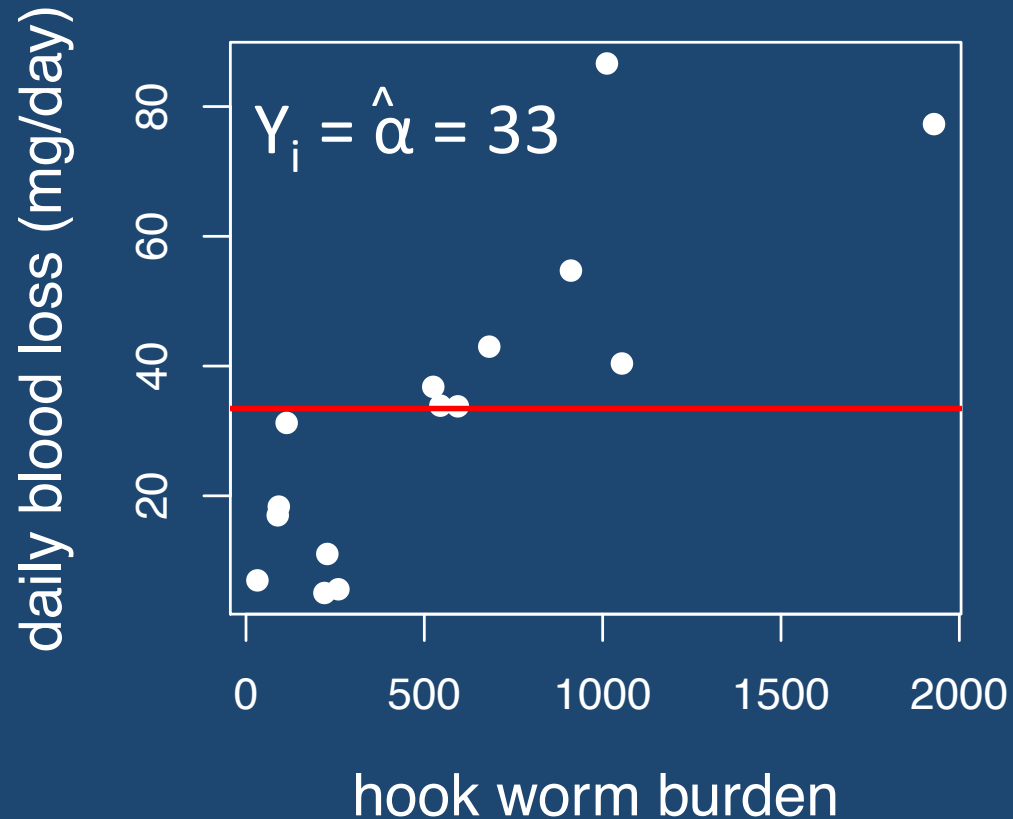How does hook worm burden affect blood loss?

Is there any relationship?



Data in Epicalc R Library taken from Areekul et al. (1970).

# Linear Regression

Null hypothesis: No relationship

$$Y = \alpha$$

Is this a good fit?

How can we get a better fit, or the best fit?



$$Y_i = \hat{\alpha} = 33$$

daily blood loss (mg/day)

hook worm burden

# Linear Regression

Null hypothesis: No relationship

$$Y_i = \alpha + \varepsilon_i$$

Is this a good fit?

How can we get a better fit, or the best fit?



$$Y_i = \hat{\alpha} = 33$$

residuals

Model

daily blood loss (mg/day)

hook worm burden

One option is Least Squares Fitting

Choose a line $Y = \hat{\alpha} + \hat{\beta}X$ to minimize $\Sigma(\text{residuals})^2$
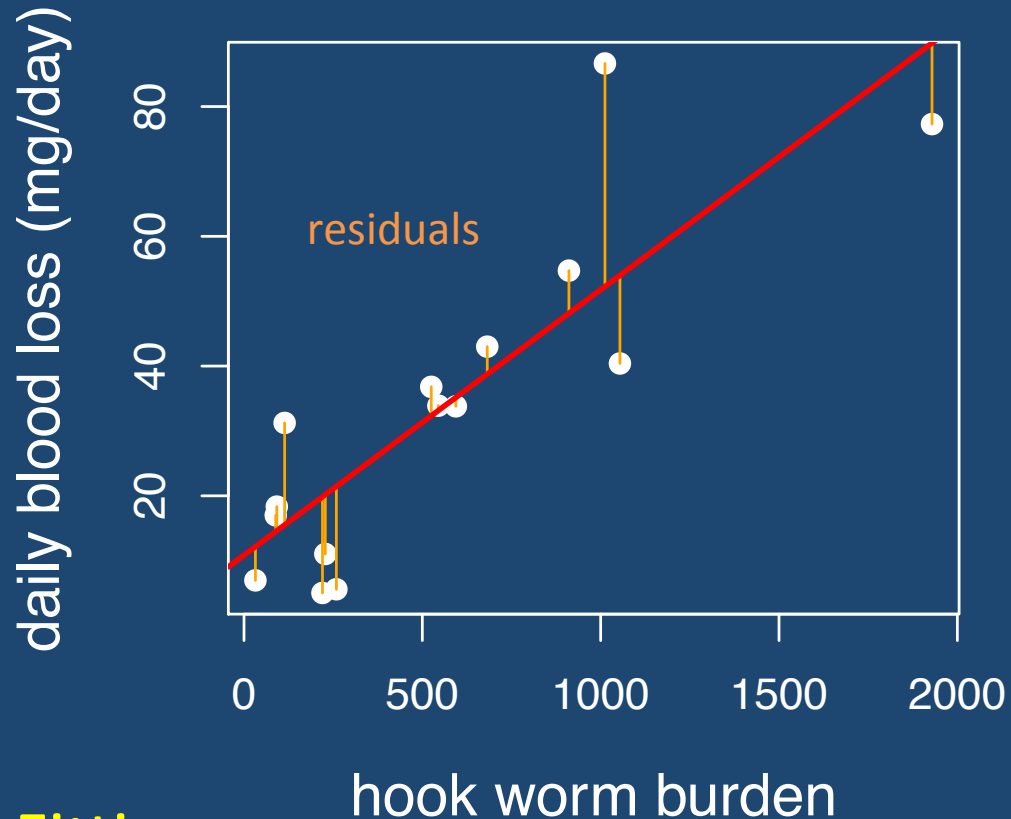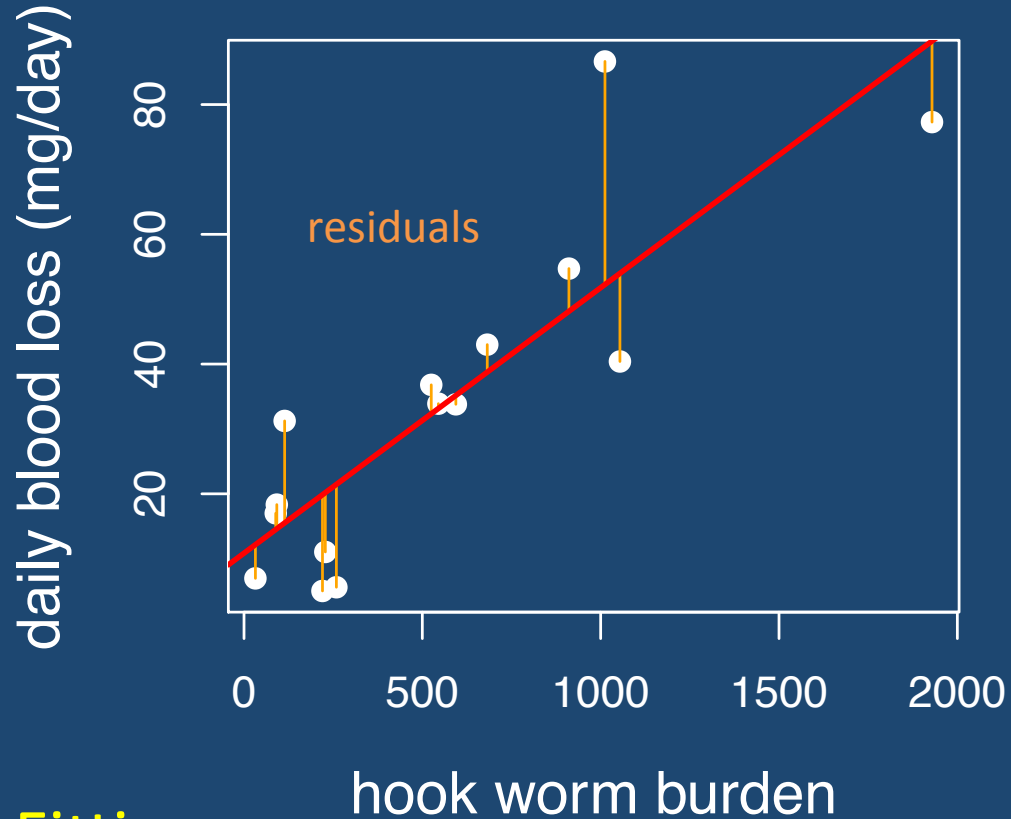
23
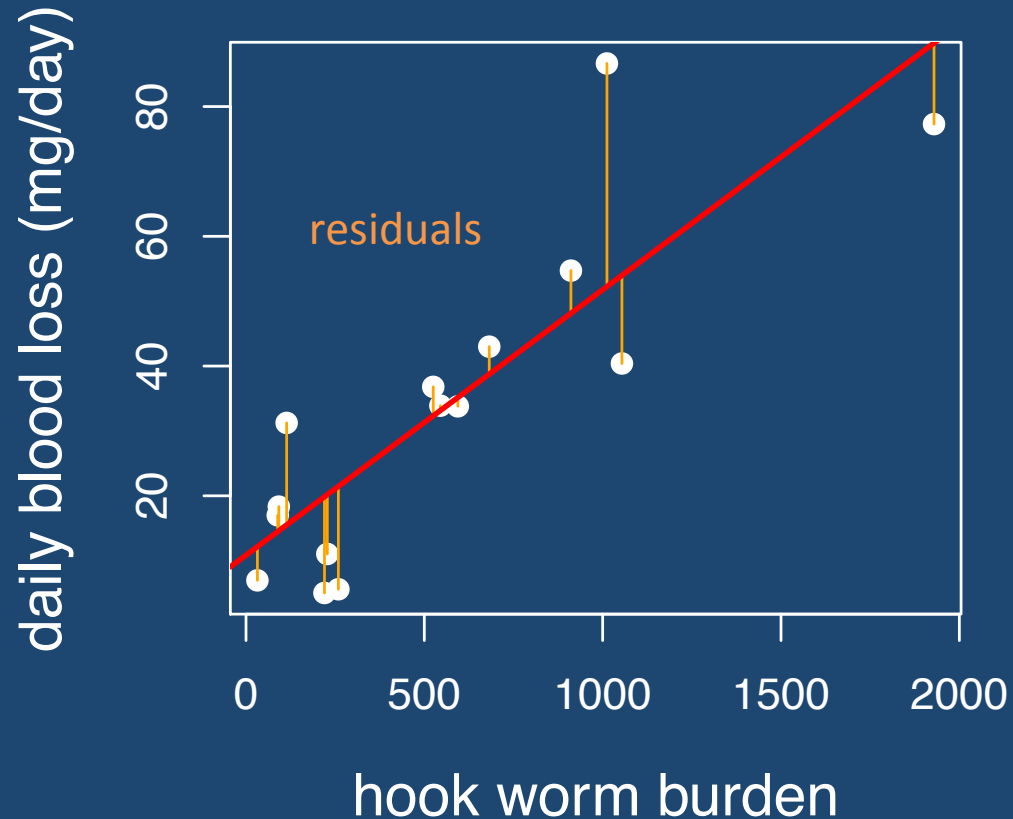
# Linear Regression

Null hypothesis: No relationship

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

Is this a good fit?

How can we get a better fit, or the best fit?

One option is Least Squares Fitting

Choose a line $Y = \hat{\alpha} + \hat{\beta}X$ to minimize $\Sigma(\text{residuals})^2$

residuals

daily blood loss (mg/day)

hook worm burden

# Linear Regression

expected daily blood loss

hook worm burden

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

intercept

effect of hook worm burden

error



daily blood loss (mg/day)

residuals

hook worm burden

One option is Least Squares Fitting

Choose a line $Y = \hat{\alpha} + \hat{\beta}X$ to minimize $\Sigma(\varepsilon_i)^2$

# Linear Regression

Another option is

Maximum Likelihood

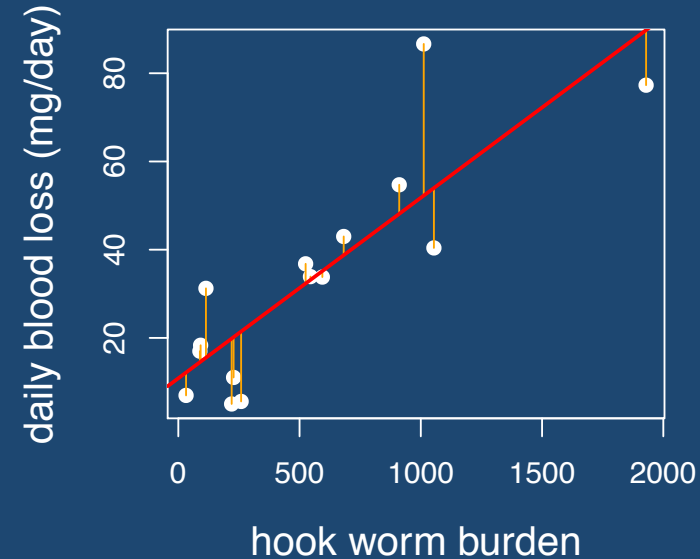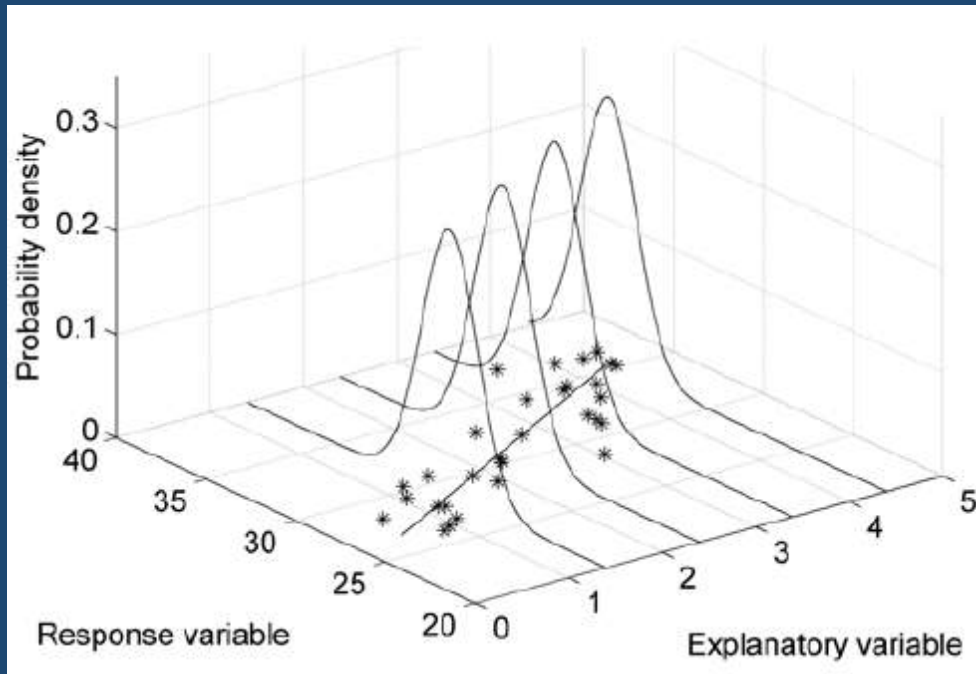$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma^2)$$
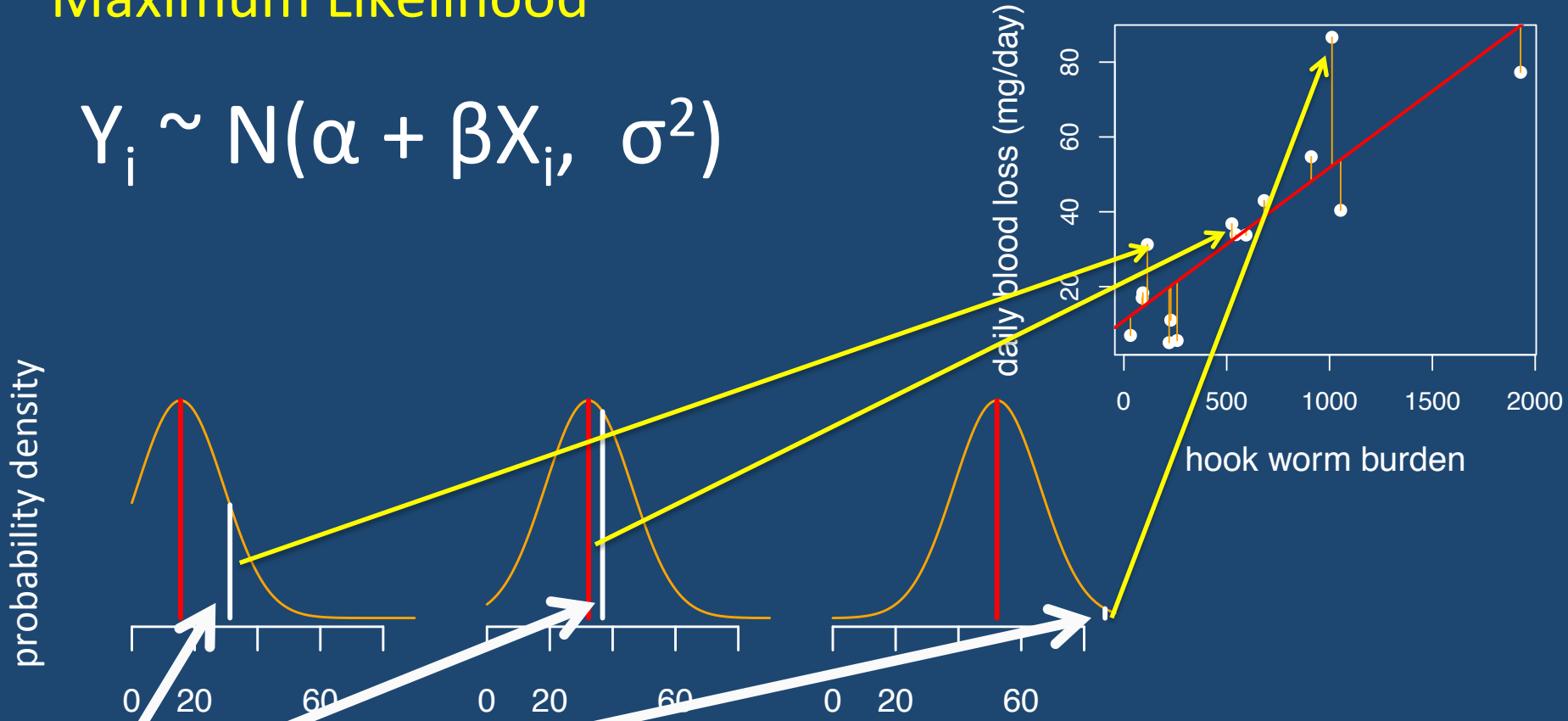


daily blood loss (mg/day) vs hook worm burden, residuals

Choose $\hat{\alpha}$, $\hat{\beta}$, $\hat{\sigma}$ to maximize the likelihood
   i.e. probability of observed data given a model

# Linear Regression

## Maximum Likelihood
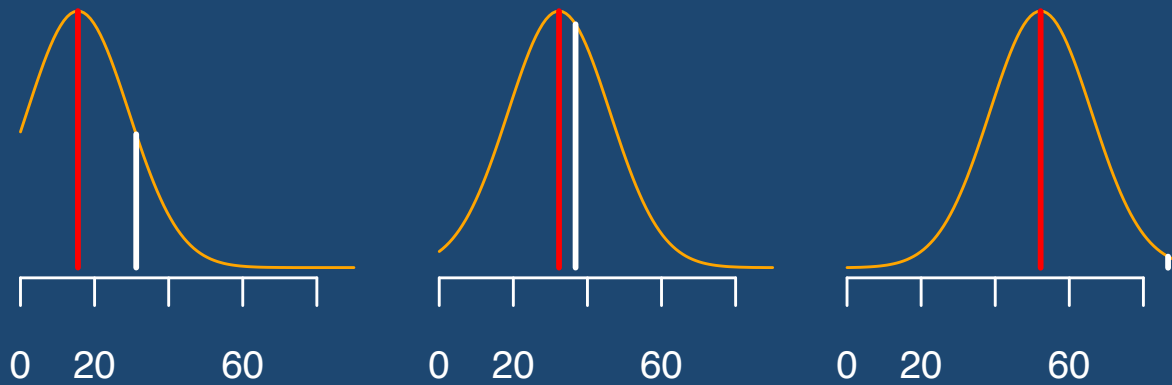
$$Y_i \sim N(\alpha + \beta X_i, \ \sigma^2)$$





Choose $\hat{\alpha}$, $\hat{\beta}$, $\hat{\sigma}$ to maximize the likelihood

i.e. probability of observed data given a model

# Linear Regression

$$Y_i \sim N(\alpha + \beta X_i, \ \sigma^2)$$

daily blood loss (mg/day)

hook worm burden

probability density

0  20  60    0  20  60    0  20  60

$$P(Y_i \mid \hat{\alpha}, \hat{\beta}, \hat{\sigma}) = \frac{1}{\hat{\sigma}\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{Y_i - (\hat{\alpha} + \hat{\beta}X_i)}{\hat{\sigma}}\right)^2}$$

28

# Linear Regression

## Maximum Likelihood

$$Y_i \sim N(\alpha + \beta X_i, \ \sigma^2)$$

daily blood loss (mg/day)

hook worm burden

probability density

0   20      60

0   20      60

0   20      60

$$P(Y_1, \ldots, Y_n \mid \hat{\alpha}, \hat{\beta}, \hat{\sigma}) = \prod_{i=1}^{n} P(Y_i \mid \hat{\alpha}, \hat{\beta}, \hat{\sigma})$$

# Linear Regression

## Maximum Likelihood



function of data

$$\downarrow$$

PDF: $\quad P(Y_1,...,Y_n \mid \hat{\alpha}, \hat{\beta}, \hat{\sigma}) = \prod_{i=1}^{n} P(Y_i \mid \hat{\alpha}, \hat{\beta}, \hat{\sigma})$

LIKELIHOOD: $\quad L(\hat{\alpha}, \hat{\beta}, \hat{\sigma} \mid Y_1,...,Y_n) = \prod_{i=1}^{n} P(Y_i \mid \hat{\alpha}, \hat{\beta}, \hat{\sigma})$

function of parameters

# Linear Regression

## Parameter Estimation & Inference



Null hypothesis:  $\beta = 0$

$\hat{\beta} = 0.04$

P(estimating a $\beta$ this extreme | null)

   P = 6.99e-05 < 0.05,

   so we reject the null hypothesis.

### Confidence intervals

Collection of
non-rejectable null hypotheses
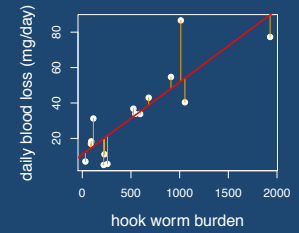
   $\hat{\beta} = 0.04$ (0.025, 0.056)

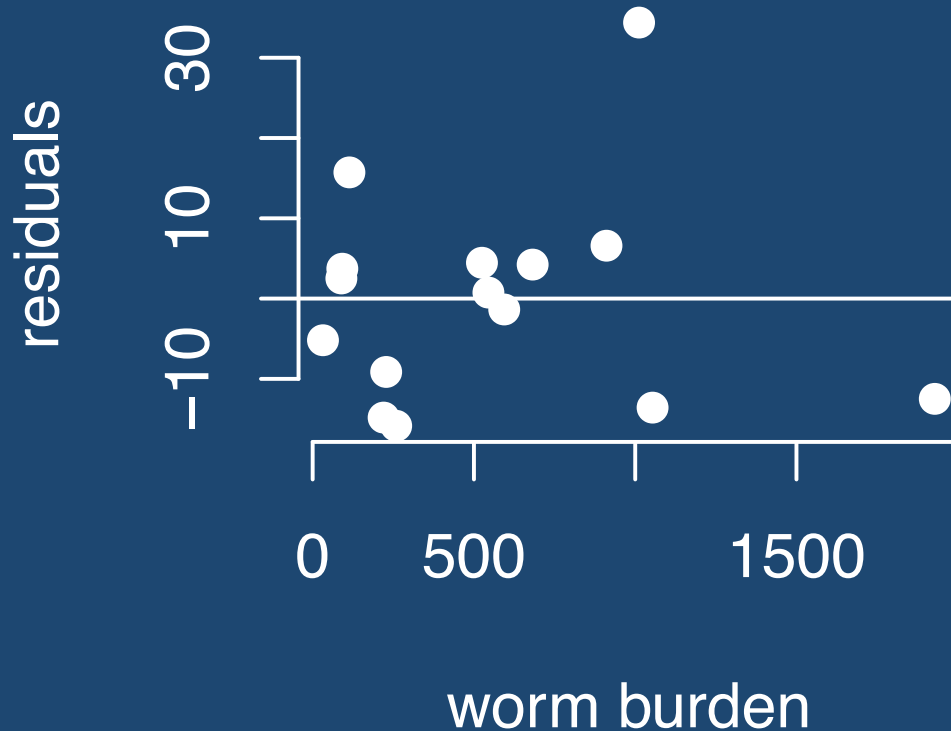# Is it a good model:
## Checking Assumptions
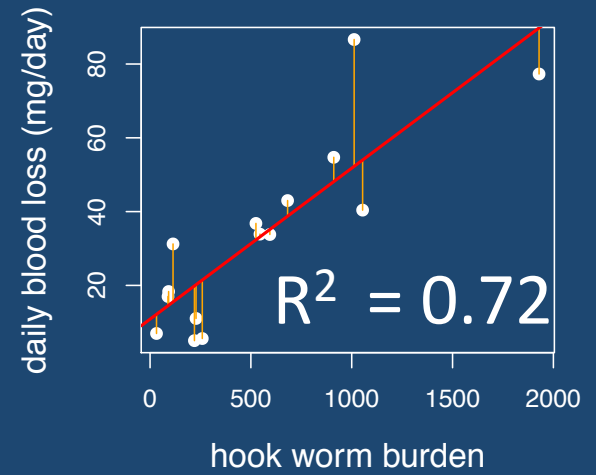


## Normality

# Is it a good model:
## Checking Assumptions



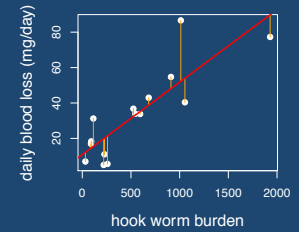Linearity        Independence        Constant Variance

# Is it a good model:
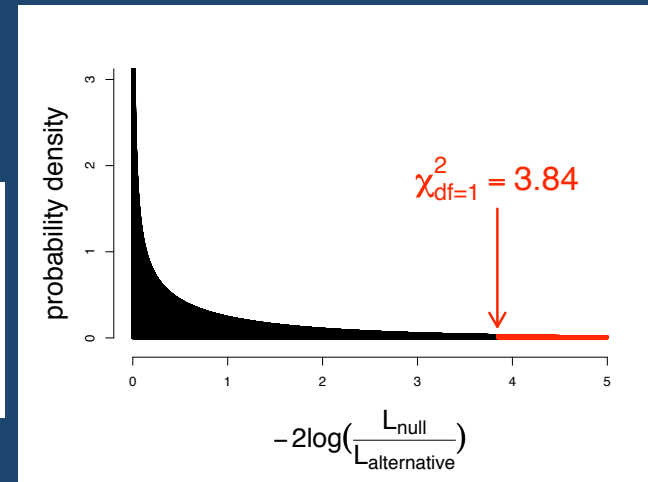# Goodness of Fit



$R^2 = $ (correlation coefficient)$^2$

How much of the variation in Y is explained by the model?
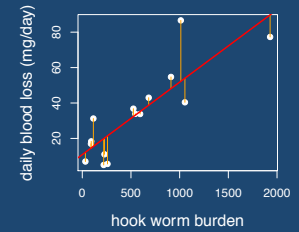
# Is it a good model:
## Goodness of Fit



## Chi Squared
## Goodness of Fit Test

$$\chi^2 = \frac{1}{n-1}\sum_{i=1}^{n}\frac{(Observed_i - Expected_i)^2}{\sigma^2}$$



- Does the observed data differ significantly from our model?
- If not, then we cannot reject our model as a bad model.
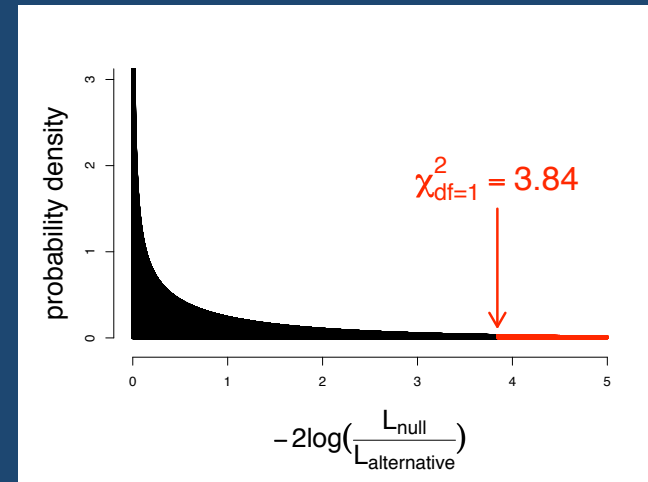- But we cannot accept our model (the null hypothesis) !
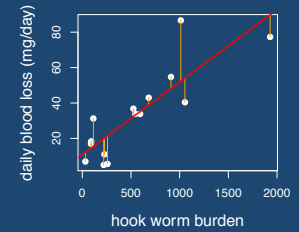
35

# Is it a good model:
# Goodness of Fit



Likelihood Ratio Test (G test, Analysis of Deviance, ANOVA)



Under the null hypothesis:

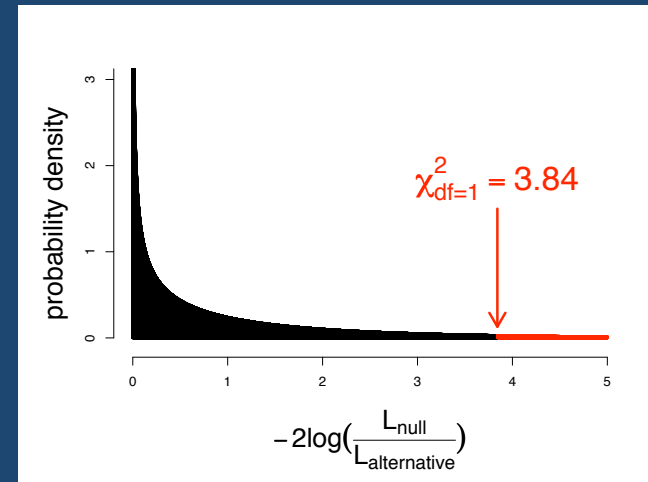$$2\log\frac{L_{MLE}}{L_{Null}} \sim \chi^2_{\text{df = difference in \# of parameters}}$$
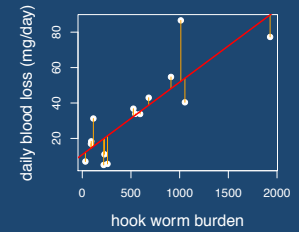
# Is it a good model:
## Model Selection



Likelihood Ratio Test (G test, Analysis of Deviance, ANOVA)



Under the null hypothesis:

$$2\log\frac{L_{\text{more parameters}}}{L_{\text{less parameters}}} \sim \chi^2_{\text{df = difference in \# of parameters}}$$

# Is it a good model:
## Model Selection

Akaike's Information Criterion (AIC)

AIC =  -2log(L) + 2(# of parameters)

penalty for adding parameters

Rank proposed models by AIC: lowest is best.

All models within 2 of lowest should be considered.

# Overfitting

- You can always fit N data points with N parameters.

- How many is too many?

- Bias/Variance Tradeoff
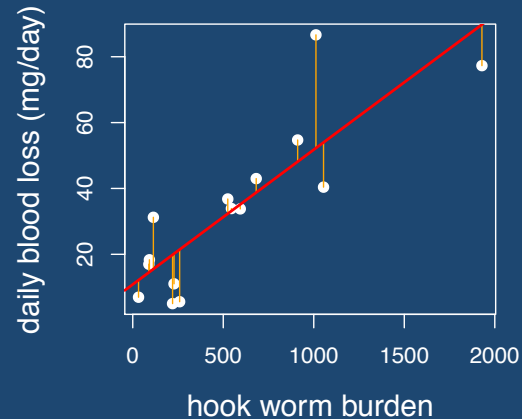
- AIC, Cross-validation

# Collinearity

- Independent variables that vary with each other

# Non-Identifiability

- Multiple parameter sets fit about equally well

# What did we just do?

- Asked a question about a relationship

- Made some observations (data)

- Formulated the relationship into a model

- Fitted the model to data

- Assessed model fit/quality (model selection)

- Inference/parameter estimation
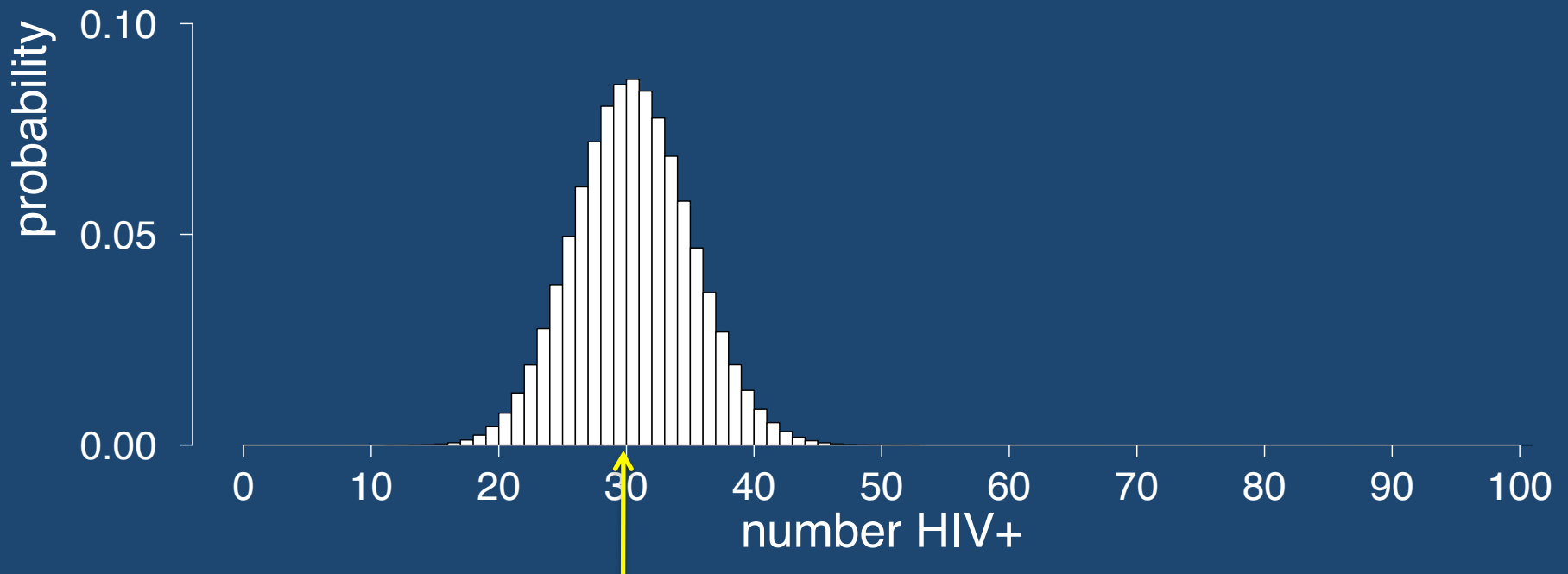
- Improved our understanding of the world





41

# Outline

1. Recap: Classical and Mechanistic Epidemiology

2. Why fit models to data?

3. Review of Linear Regression

4. Maximum Likelihood and Fitting Simple Models

5. Fitting Dynamic Models to Data

6. Summary

In a population of 1,000,000 people with a true prevalence of 30%, the probability distribution of number of positive individuals if 100 are sampled:
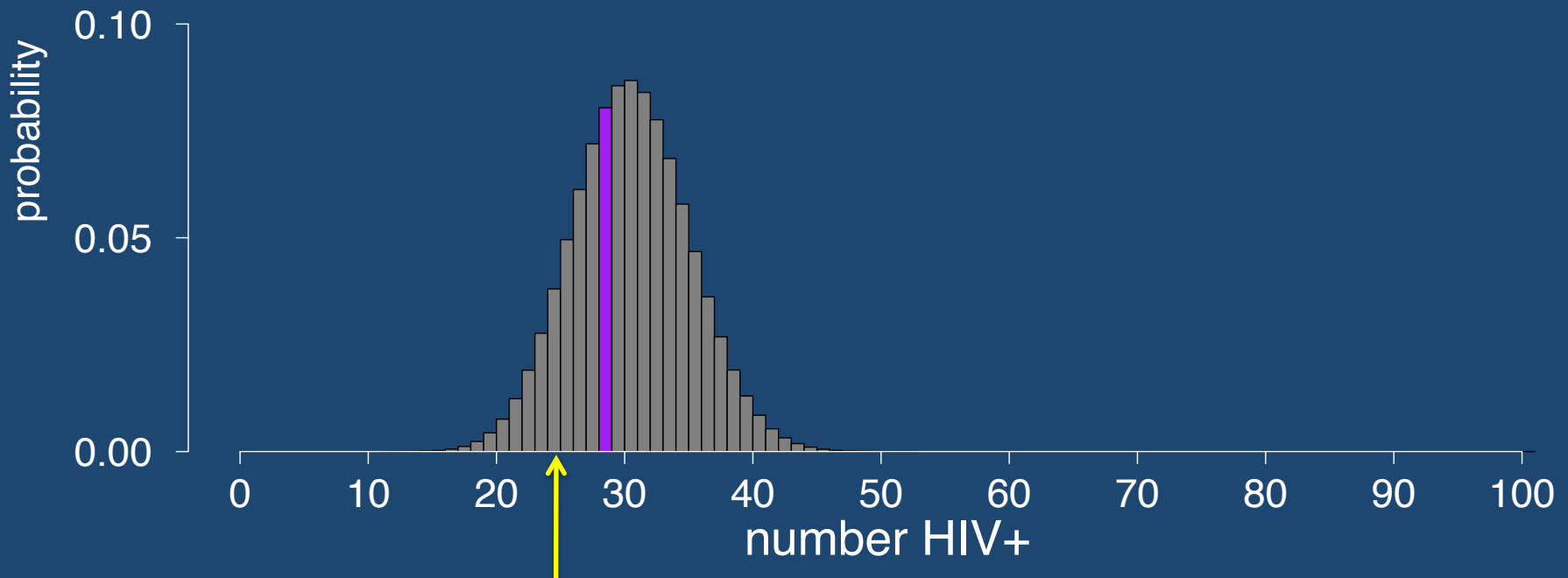
$$f(x) = \binom{100}{x}(0.3)^x(0.7)^{100-x}$$



We sample 100 people once and 28 are positive:

```
> rbinom(n = 1, size = 100, prob = .3)
[1] 28
```

# Introduction to Likelihood

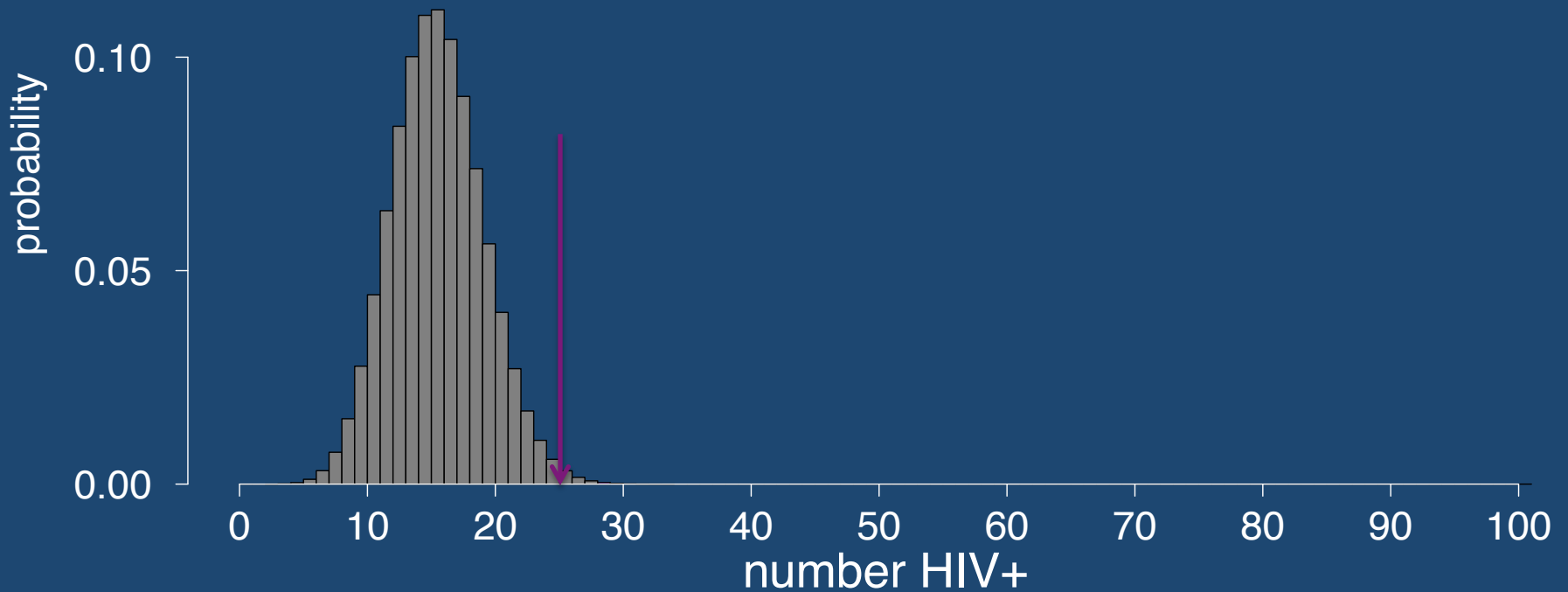**hypothetical   prevalence: 30 %**

dbinom(28, 100, 0.3) = 0.0804



We sample 100 people once and 28 are positive:

> rbinom(n = 1, size = 100, prob = .3)
[1] 28

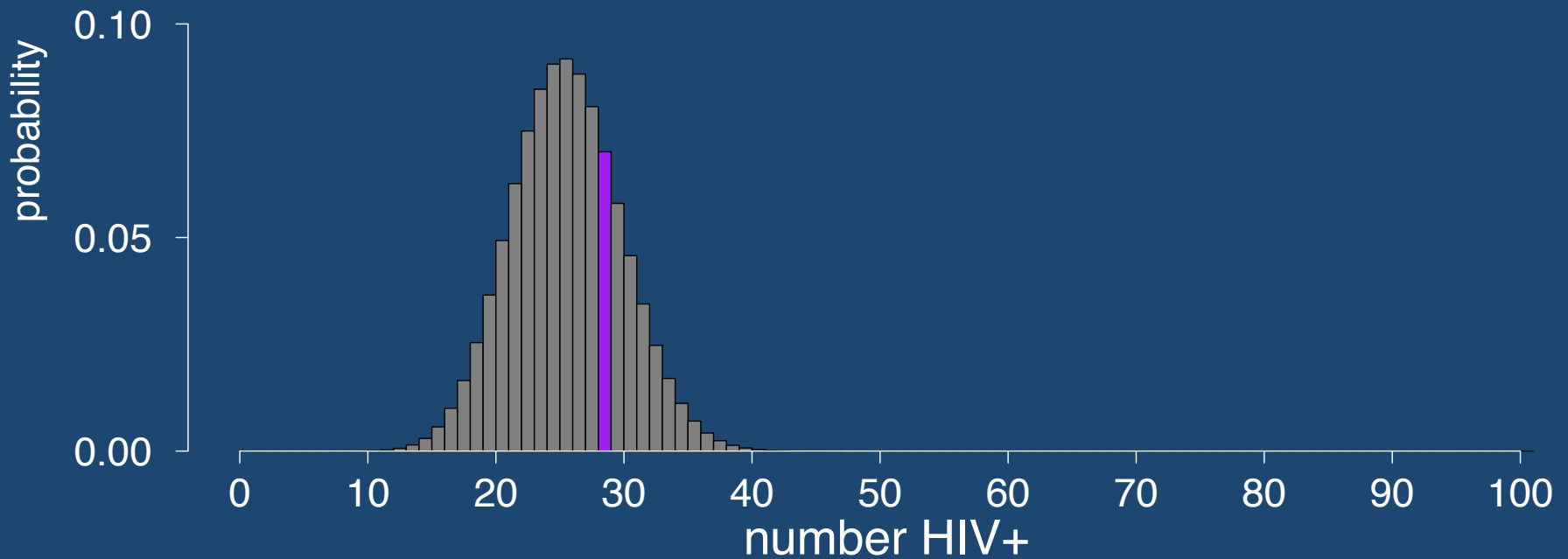**hypothetical   prevalence: 15 %**

dbinom(28, 100, 0.15) = 0.000353

45

hypothetical    prevalence: 20 %

dbinom(28, 100, 0.2) = 0.0141

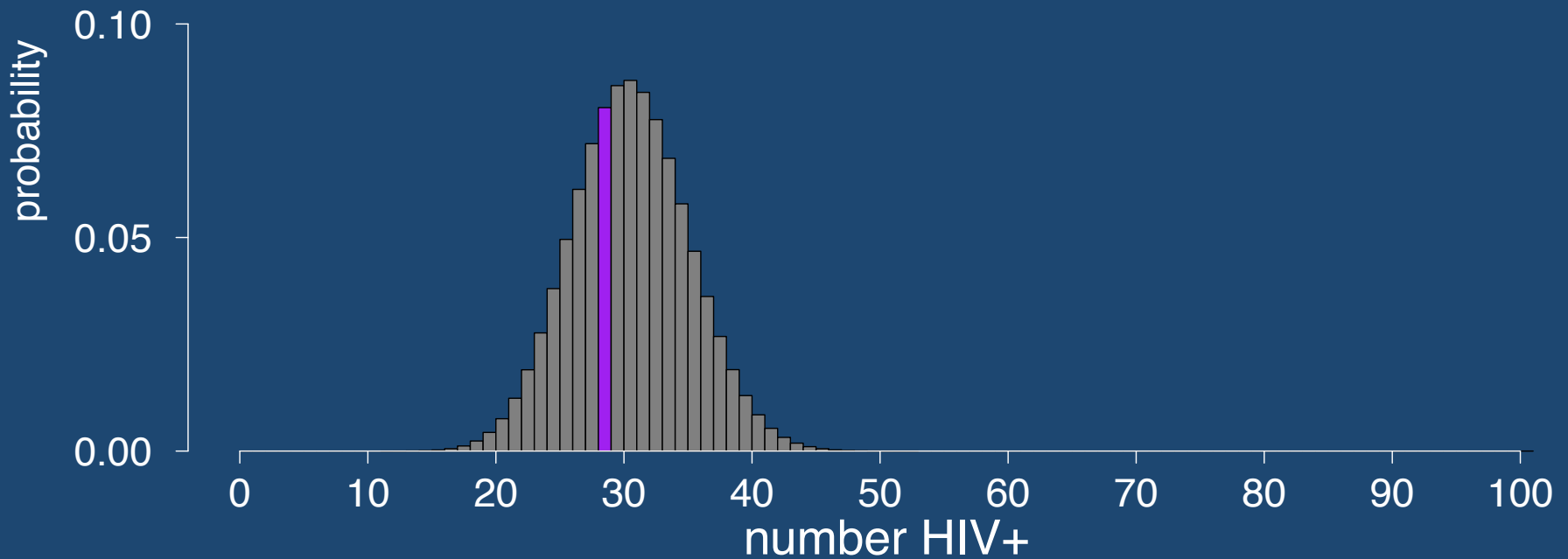hypothetical   prevalence: 25 %

dbinom(28, 100, 0.25) = 0.0701

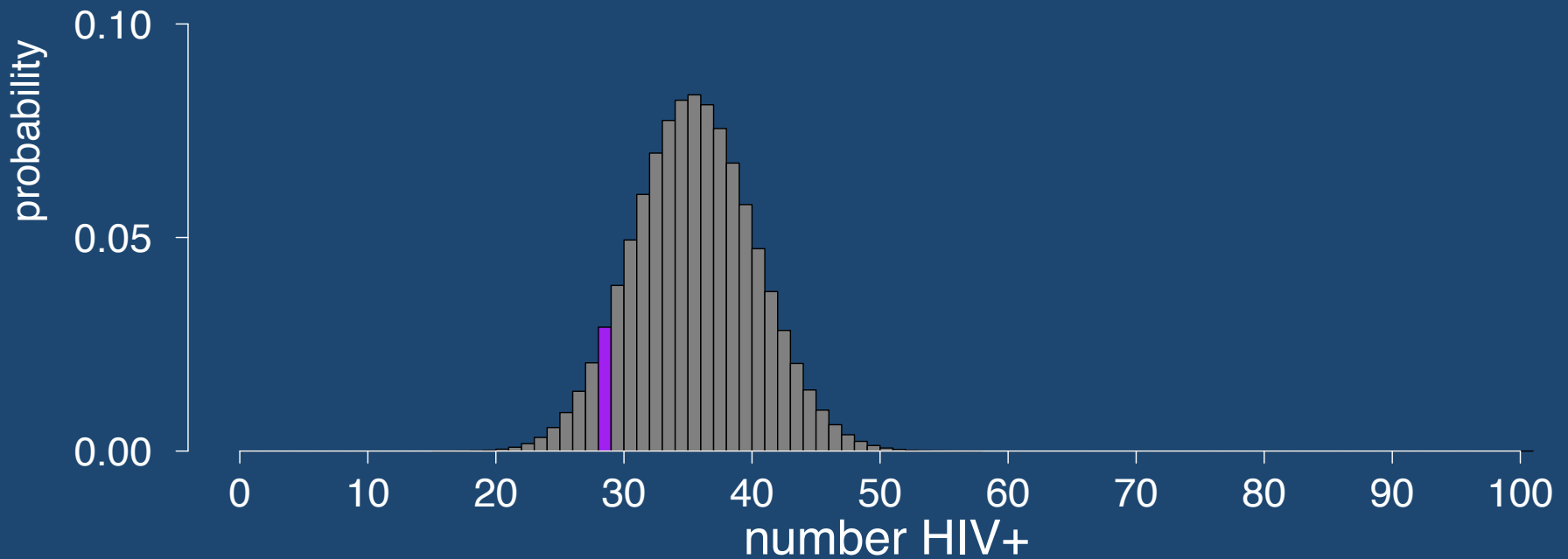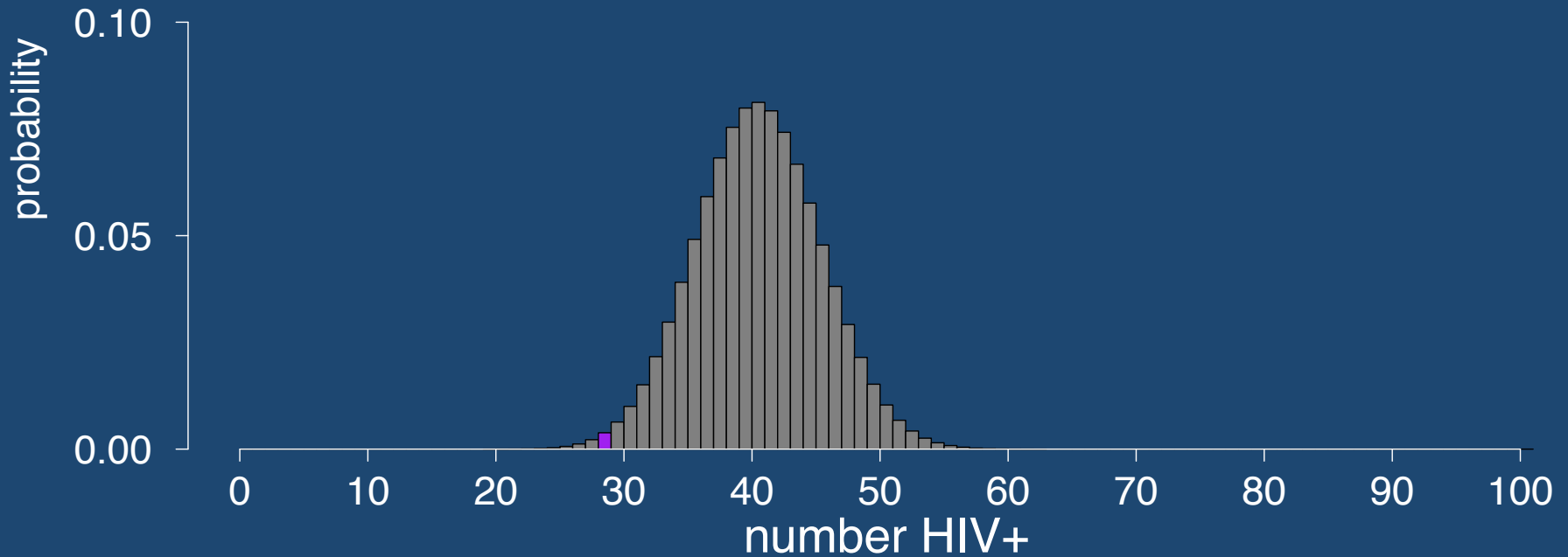hypothetical prevalence: 30 %

dbinom(28, 100, 0.3) = 0.0804

**hypothetical   prevalence: 35 %**

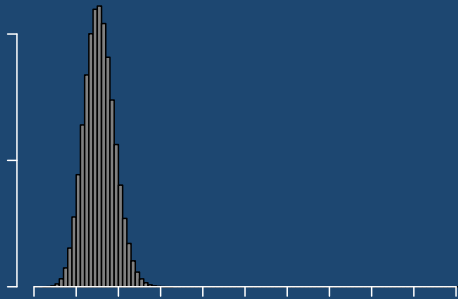dbinom(28, 100, 0.35) = 0.029

**hypothetical prevalence: 40 %**

dbinom(28, 100, 0.4) = 0.00383

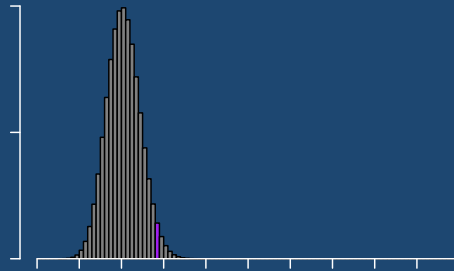# Which prevalence gives the greatest probability of observing exactly 28/100?



probability

number HIV+

51

# Which of these prevalence values is most likely given our data?

p(our data given prevalence) = LIKELIHOOD

Maximum Likelihood Estimate
parameter value giving greatest probability
of the data having occurred.

What do you think is the MLE here?

MLE = 28/100 = 0.28

hypothetical prevalence

true unknown value = 0.30

different null hypotheses

53

# Defining Likelihood

- L(parameter | data) = p(data | parameter)

- Not a probability distribution.

- Probabilities taken from many different distributions.

function of x

$$\text{PDF:} \quad f(x|p) = \binom{n}{x} (p)^x (1-p)^{n-x}$$

$$\text{LIKELIHOOD:} \quad L(p|x) = \binom{n}{x} (p)^x (1-p)^{n-x}$$

function of p

# Deriving the Maximum Likelihood Estimate

maximize

$$L(p) = \binom{n}{x}(p)^x(1-p)^{n-x}$$

maximize

$$\log\big(L(p)\big) = \log\left[\binom{n}{x}(p)^x(1-p)^{n-x}\right]$$

minimize

$$l(p) = -\log\left[\binom{n}{x}(p)^x(1-p)^{n-x}\right]$$

likelihood / hypothetical prevalence

log(likelihood) / hypothetical prevalence

−log(likelihood) / hypothetical prevalence

55

# Likelihood



Maximum Likelihood Estimate

$$\hat{p} = \frac{x}{n} = \frac{28}{100} = 0.28$$

hypothetical prevalence

we usually minimize the –log(likelihood)

Maximum Likelihood Estimate

$$\hat{p} = \frac{x}{n} = \frac{28}{100} = 0.28$$

hypothetical prevalence

# Building Confidence Intervals
## Likelihood Ratio Test

If the null hypothesis were true then

$$2 \log \left( \frac{L(\text{alternative hypothesis})}{L(\text{null hypothesis})} \right) \sim \chi^2_{df=1}$$

PDF for $\chi^2_{df=1}$



Why does this work?

- Adding irrelevant parameters *always* improves the fit.

- How much should fit improve due to chance alone by adding an irrelevant parameter?

- Fit improvement, as measured above, is approximately $\chi^2_{df}$ distributed with df = to the difference in parameters used to fit.

# Building Confidence Intervals
## Likelihood Ratio Test

If the null hypothesis were true then

$$2\log\left(\frac{L(\text{alternative hypothesis})}{L(\text{null hypothesis})}\right) \sim \chi^2_{df=1}$$
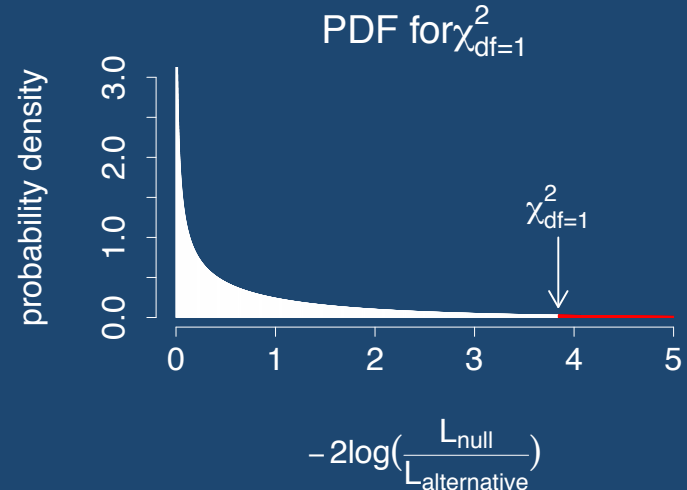
$$2\log(L_{\text{MLE}}) - 2\log(L_{\text{null}}) \sim \chi^2_{df=1}$$

$$-2l_{MLE} + 2l_{null} \sim \chi^2_{df=1}$$

PDF for $\chi^2_{df=1}$

probability density

$-2\log(\frac{L_{\text{null}}}{L_{\text{alternative}}})$
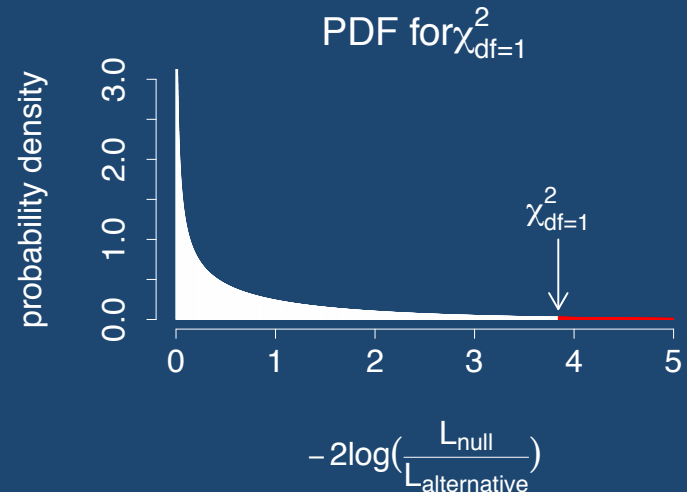
So if our α = .05, then we reject any null hypothesis for which

$$-2l_{MLE} + 2l_{null} > \chi^2_{df=1,\alpha=.05} = 3.84$$

> qchisq(p = .95, df = 1)
[1] 3.841459

$$l_{null} - l_{MLE} > 1.92$$

If $\log(L_{MLE}) - \log(L_{null}) > 1.92$,

we reject that null hypothesis. 59

# Building Confidence Intervals
## Likelihood Ratio Test

we usually minimize the –log(likelihood)



Maximum Likelihood Estimate

$$\hat{p} = \frac{x}{n} = \frac{28}{100} = 0.28$$

Let's zoom in…

–log(likelihood)

400

300

200

100

0

0.0    0.2    0.4    0.6    0.8    1.0

hypothetical prevalence

# Building Confidence Intervals
## Likelihood Ratio Test



Maximum Likelihood Estimate

$$\hat{p} = \frac{x}{n} = \frac{28}{100} = 0.28$$

−log(likelihood)

hypothetical prevalence

61

# Building Confidence Intervals
## Likelihood Ratio Test



$$l_{null} - l_{MLE} > 1.92$$

hypothetical prevalence

# Building Confidence Intervals
## Likelihood Ratio Test



95% CI includes HIV prevalences of 19.9% to 37.2%

# Outline

1. Recap: Classical and Mechanistic Epidemiology

2. Why fit models to data?

3. Review of Linear Regression

4. Maximum Likelihood and Fitting Simple Models

5. Fitting Dynamic Models to Data

6. Summary

## Statistical Models & Dynamic Models

**Statistical Models**

- Account for bias and random error to find correlations that may imply causality.

- Often the first step to assessing relationships.

- Assume independence of individuals (at some scale).

**Dynamic Models**

- Systems Approach: Explicitly model multiple mechanisms to understand their interactions.

- Links observed relationships at different scales.

- Explicitly focuses on dependence of individuals

By developing dynamic models in a probabilistic framework we can account for dependence, random error, and bias while linking patterns at multiple scales.

# Fitting Dynamic Models to Data

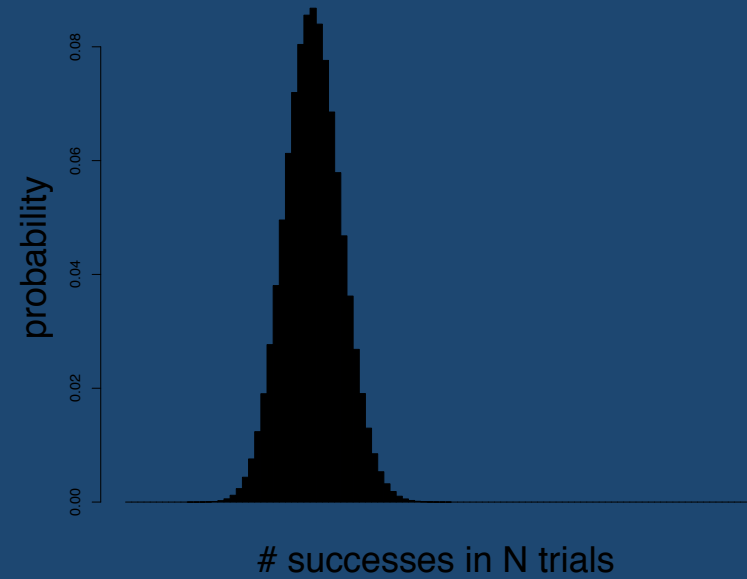Adapt our dynamic models in a probabilistic framework so we can ask:

What is the probability that a model would have generated the observed data?

What is the likelihood of a model given the data?
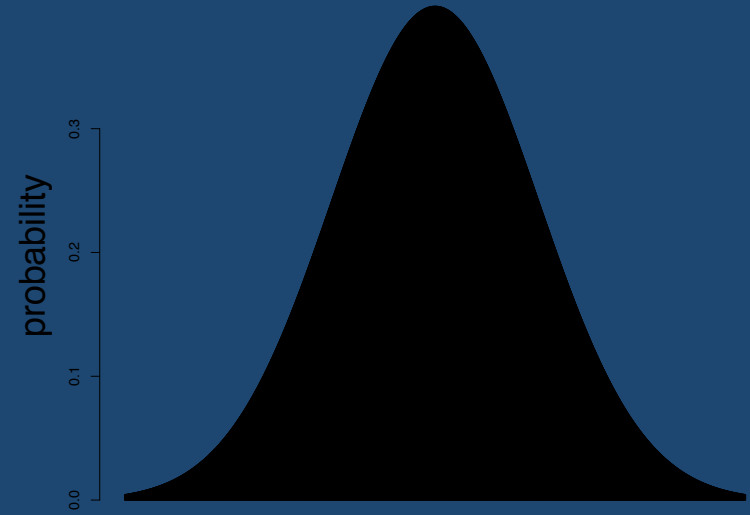
Likelihood of parameters
(given data)

67

# Binomial Distribution

probability

# successes in N trials

---

**Distribution**

↓

**Likelihood of parameters
(given data)**

68

# Normal Distribution



probability

0.3

0.2

0.1

0.0

(approximately) continuous variable

**Distribution**

**Likelihood** of parameters
(given data)

# Exponential Distribution



probability

time until event

**Distribution**

**Likelihood** of parameters
(given data)

# Poisson Distribution
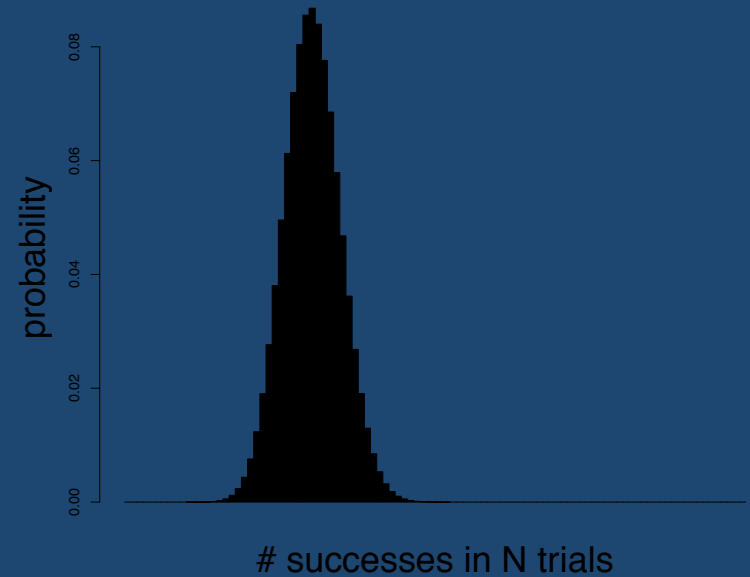


probability

# of events in time interval

**Distribution**

**Likelihood** of parameters
(given data)

# Binomial Distribution



**Stochastic Component of Model**

Distribution

Likelihood of parameters
(given data)

Binomial

HIV in Harare

# successes in N trials

probability

prevalence %

# successes in N trials

probability

Distribution

Data

Likelihood of parameters
(given data)

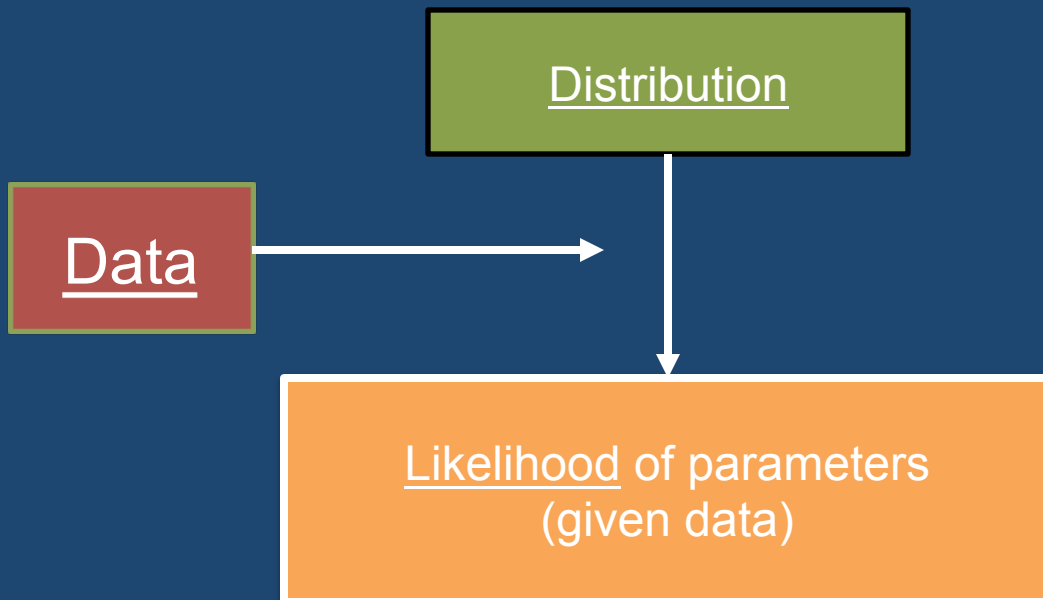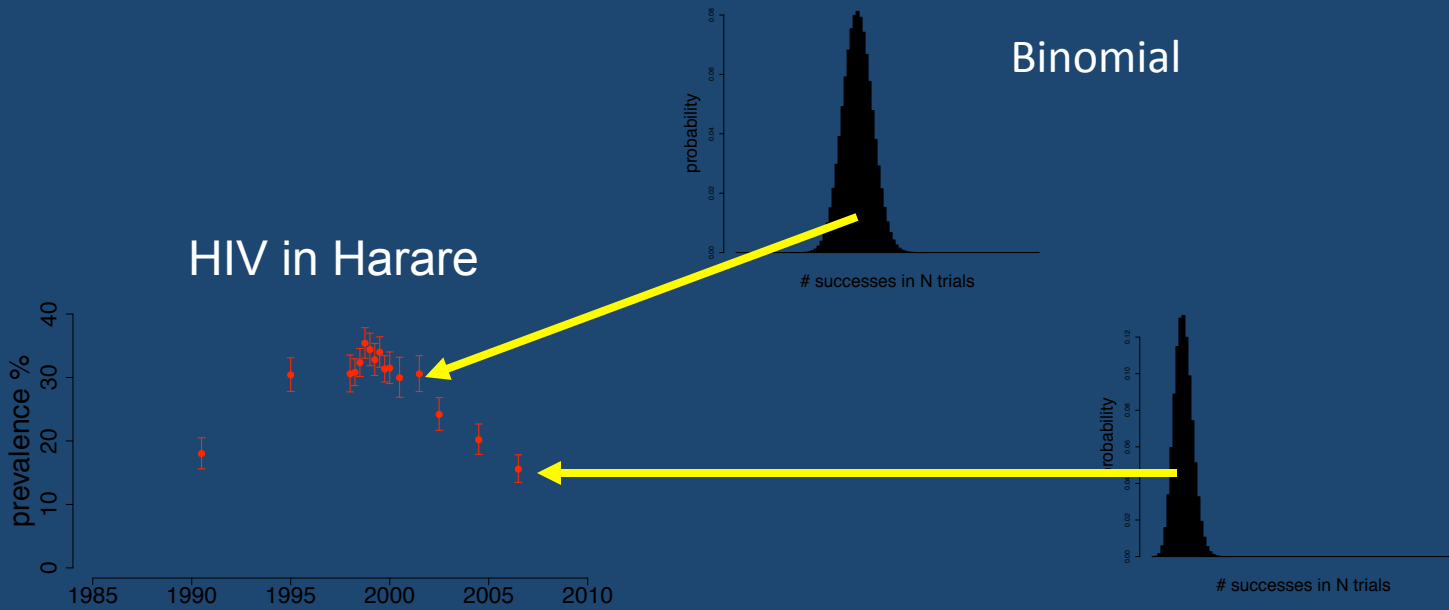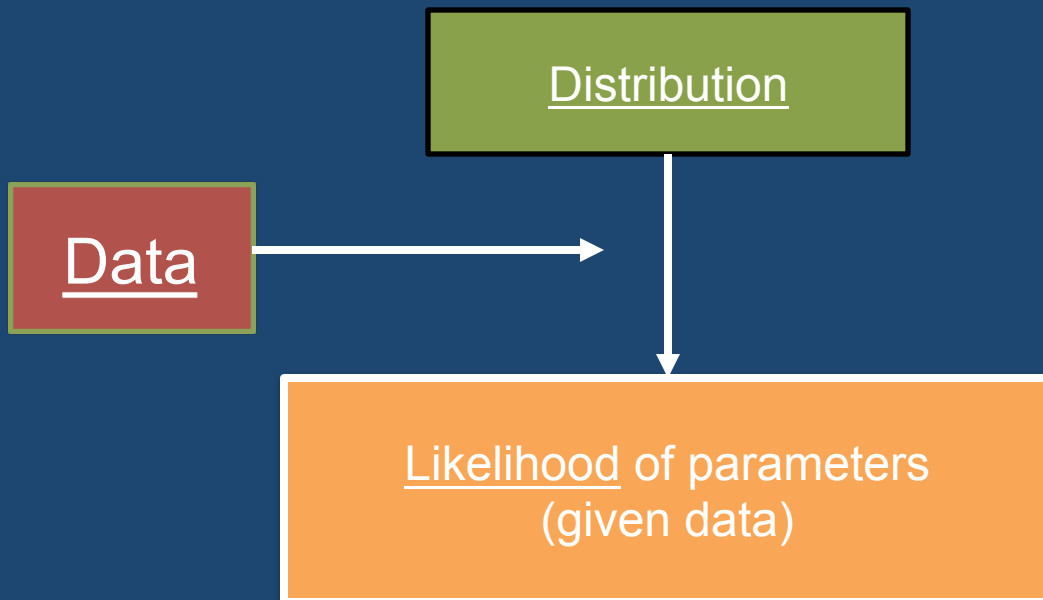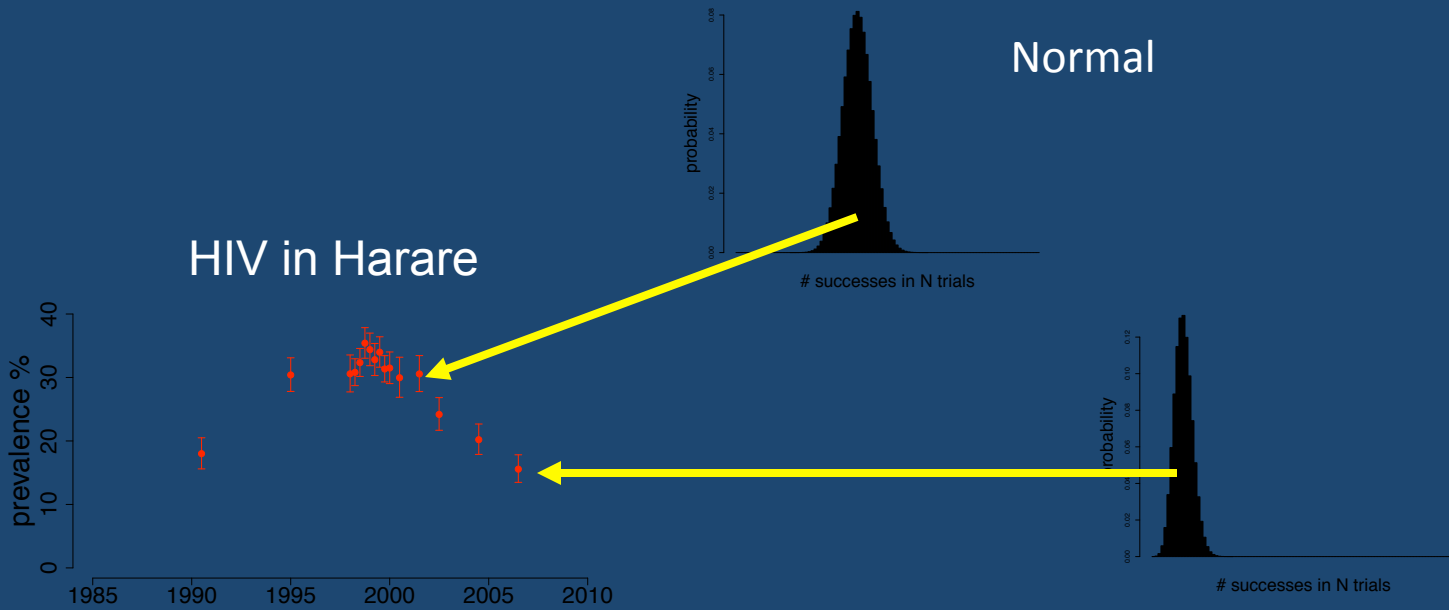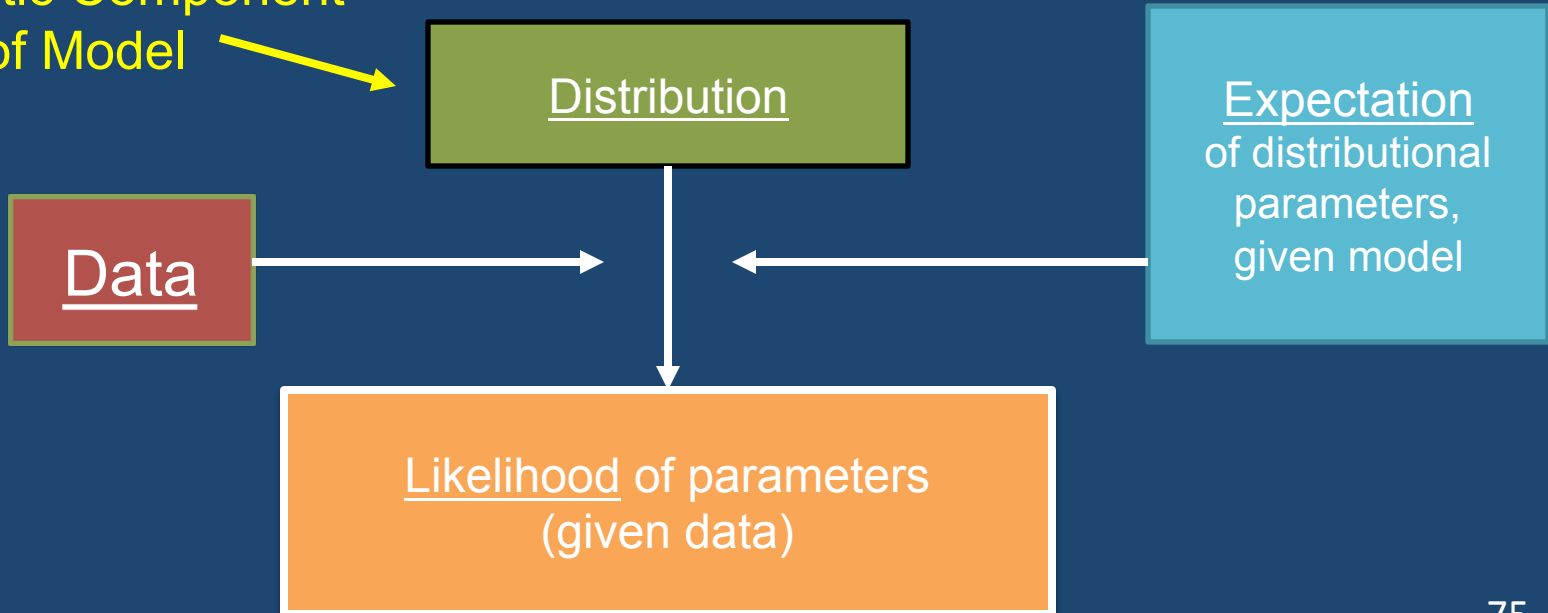Normal

HIV in Harare

Distribution

Data
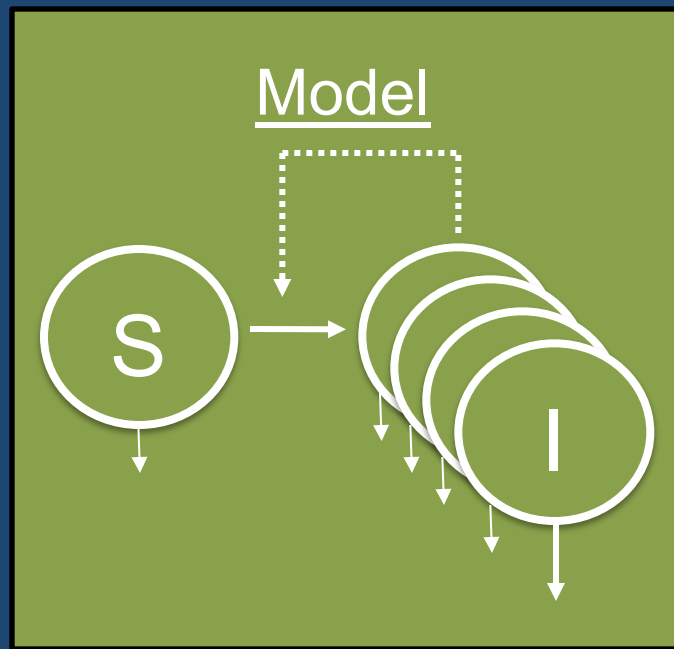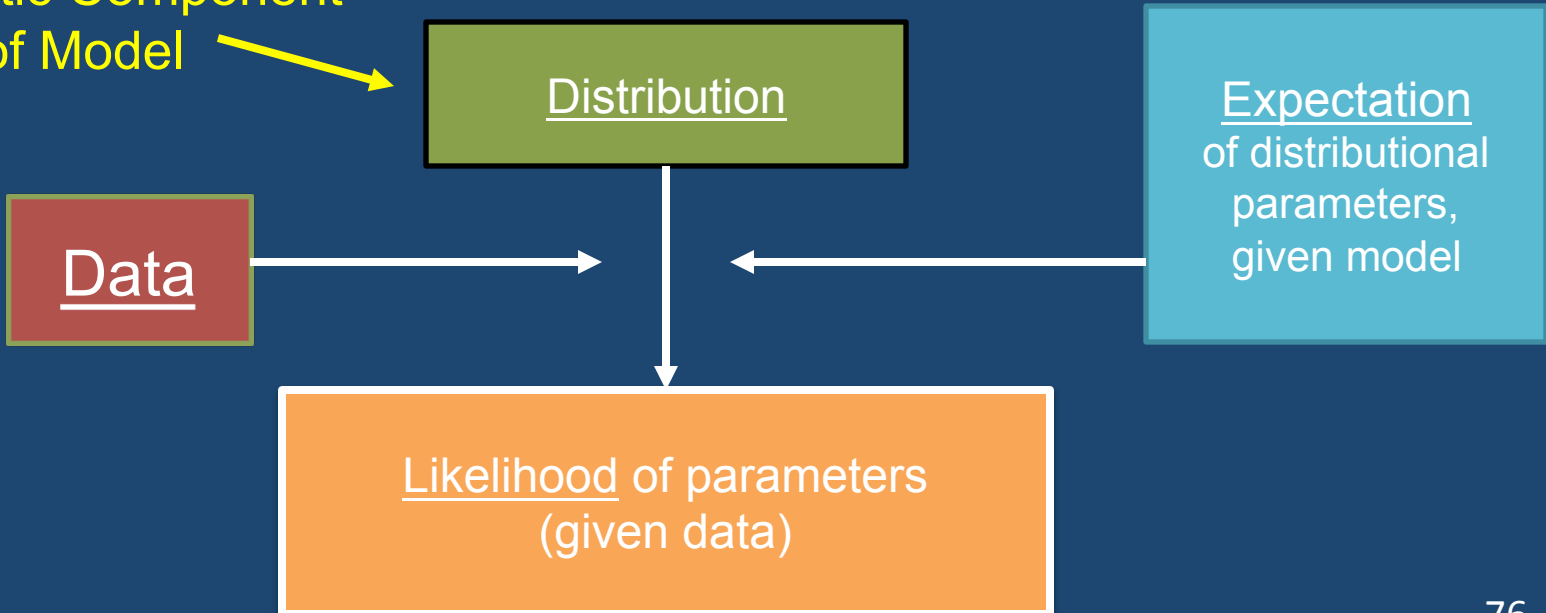
Likelihood of parameters
(given data)

74

Stochastic Component
of Model

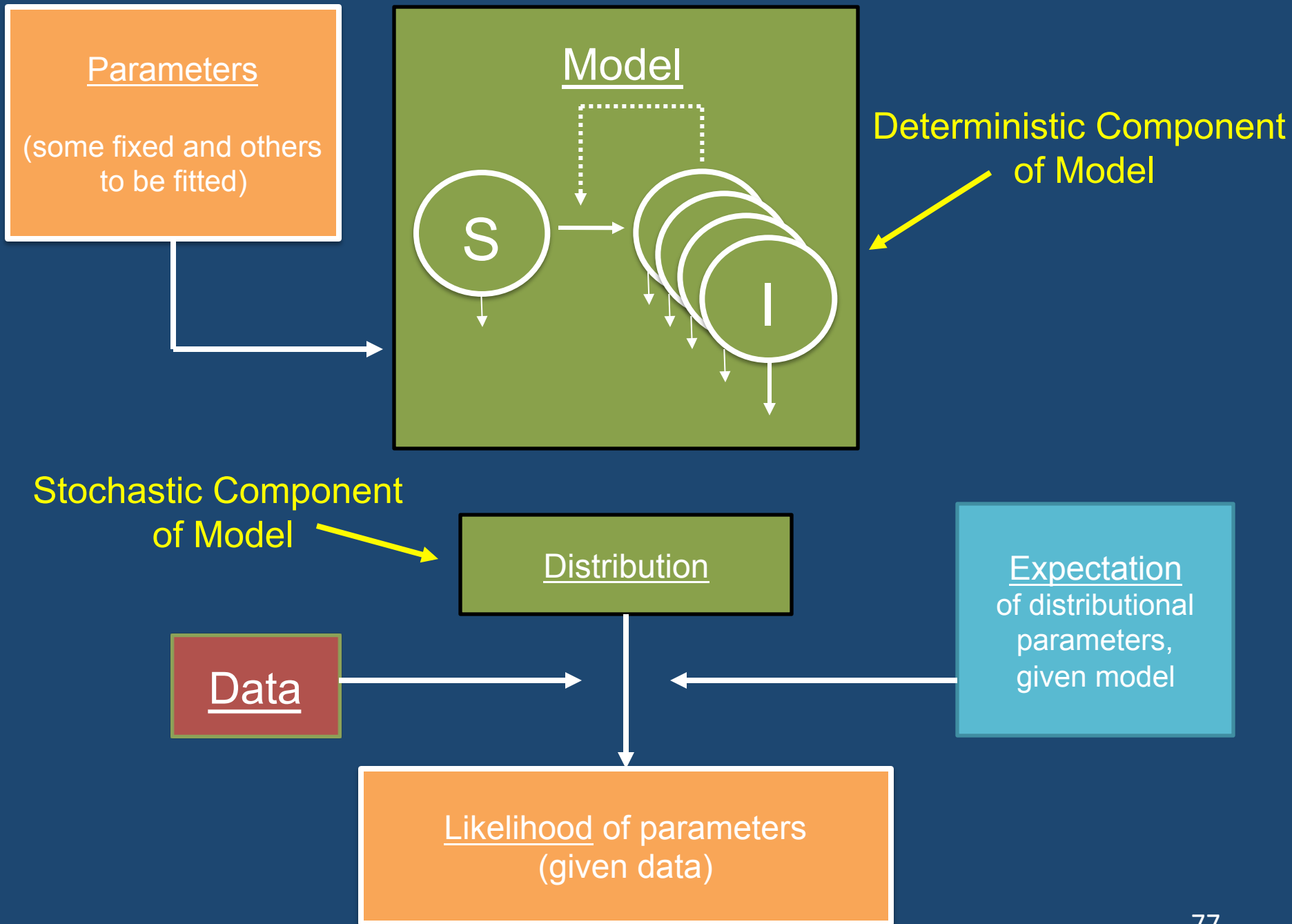Distribution

Data

Expectation
of distributional
parameters,
given model

Likelihood of parameters
(given data)

75

# Collinearity

- Independent variables that vary with each other

# Non-Identifiability

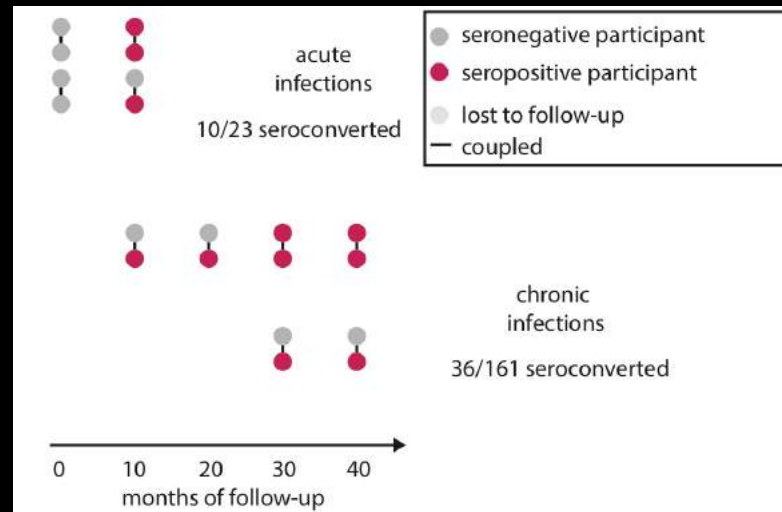- Multiple parameter sets fit about equally well

- Can be informative in dynamic models

# Rakai *Retrospective Couples* Cohort

7x as infectious for first 5 month
26x as infectious for first 3 months

$EHM_{acute}$ = 30 or 70

# Comparing Results

| Study | $RH_{acute}$ | $d_{acute}$ (months) |
|---|---|---|
| Wawer et al. (2005) | 7.25 (3.05 – 17.3) | 5 |
| Hollingsworth et al. (2008) | 26 | 2.9 (1.23-6) |

# Collinearity in Fitted Parameters



Holl. 2008: $RH_{acute} = 26$, $d_{acute} = 2.9$

y-axis: $d_{acute}$ (months)

x-axis: $RH_{acute}$

# Collinearity in Fitted Parameters



Holl. 2008:  $RH_{acute}$ = 26,   $d_{acute}$ = 2.9

our refit:  $RH_{acute}$ = 42,   $d_{acute}$ = 1.5

Refit the same model using Bayesian MCMC

# Collinearity in Fitted Parameters
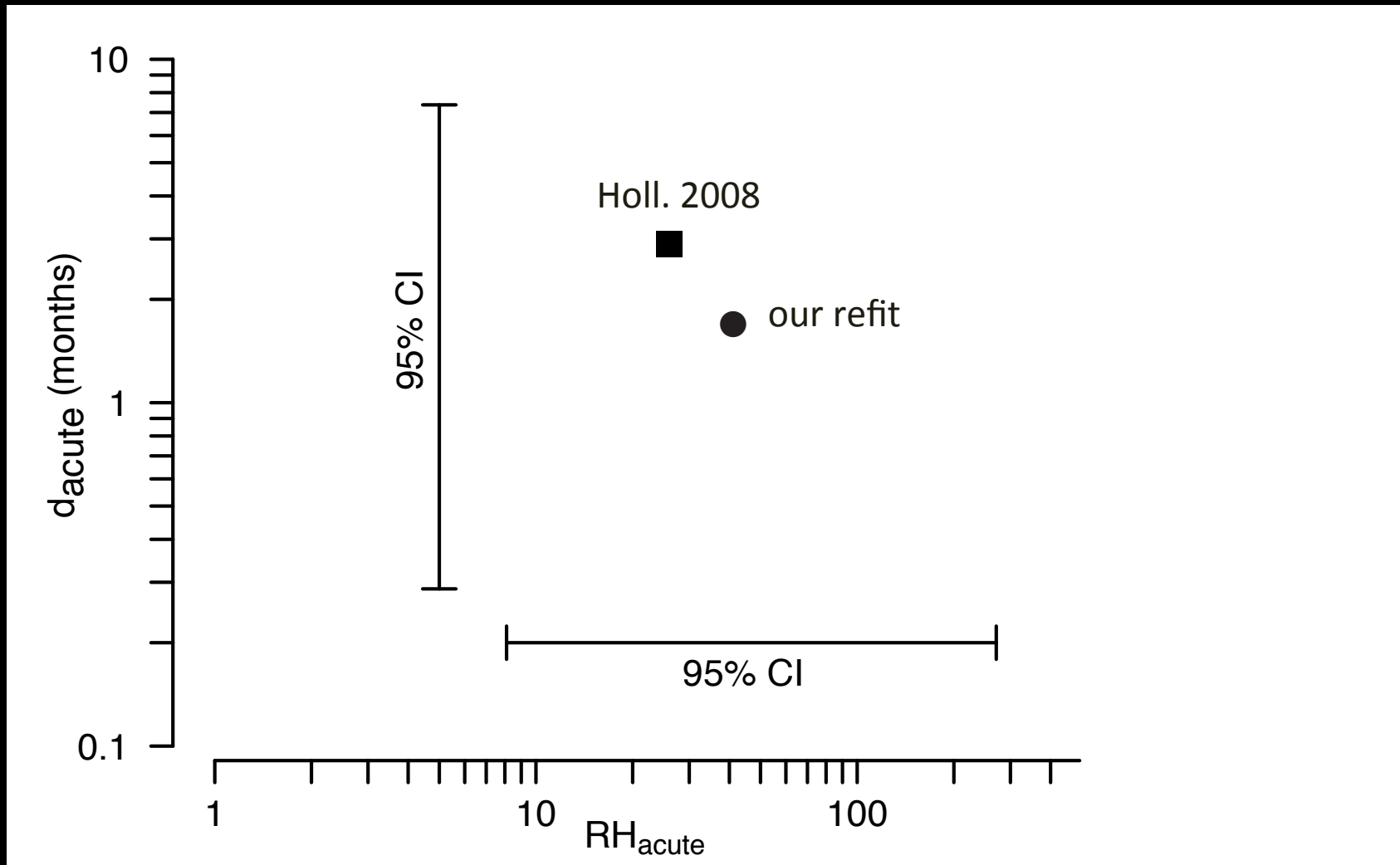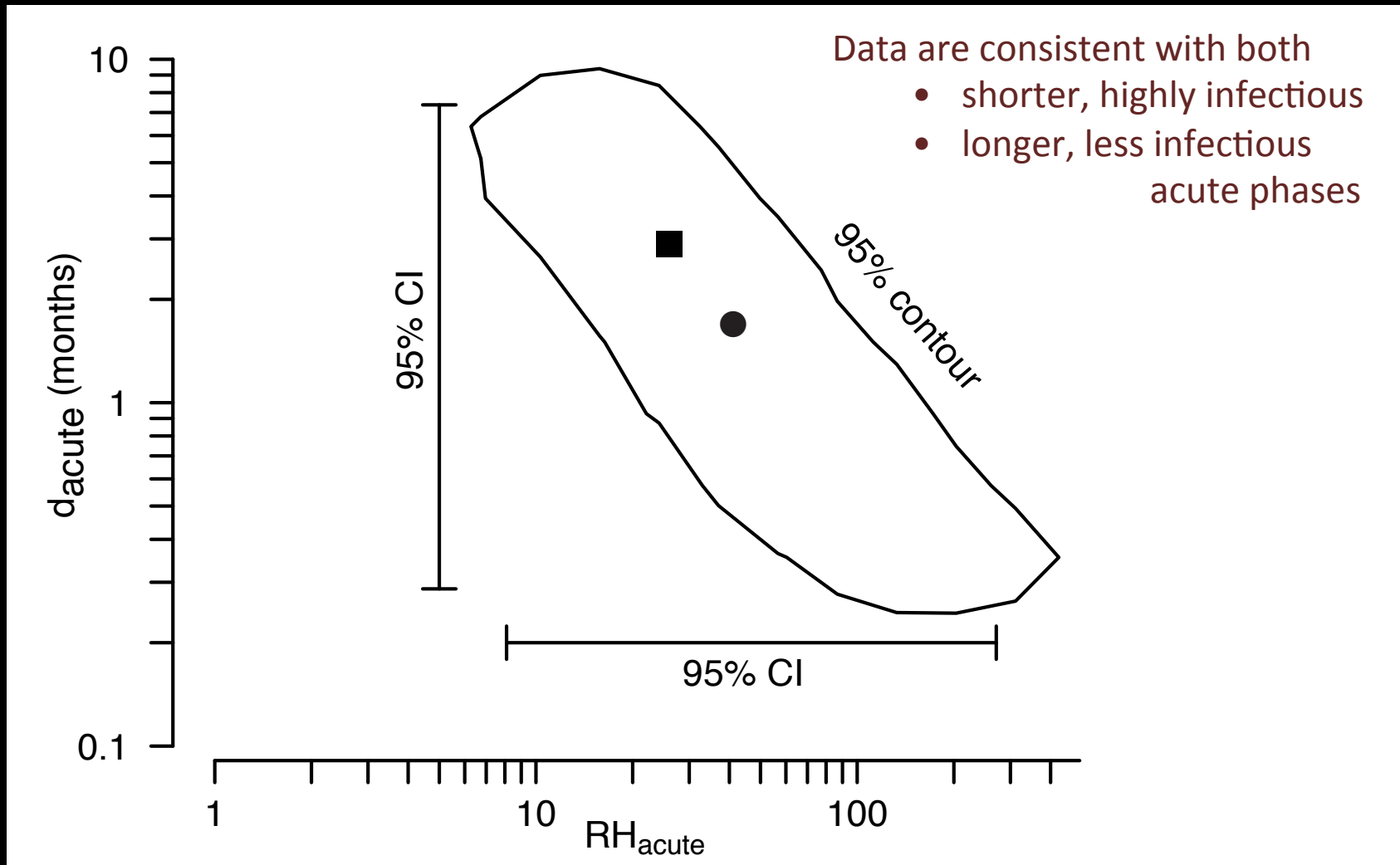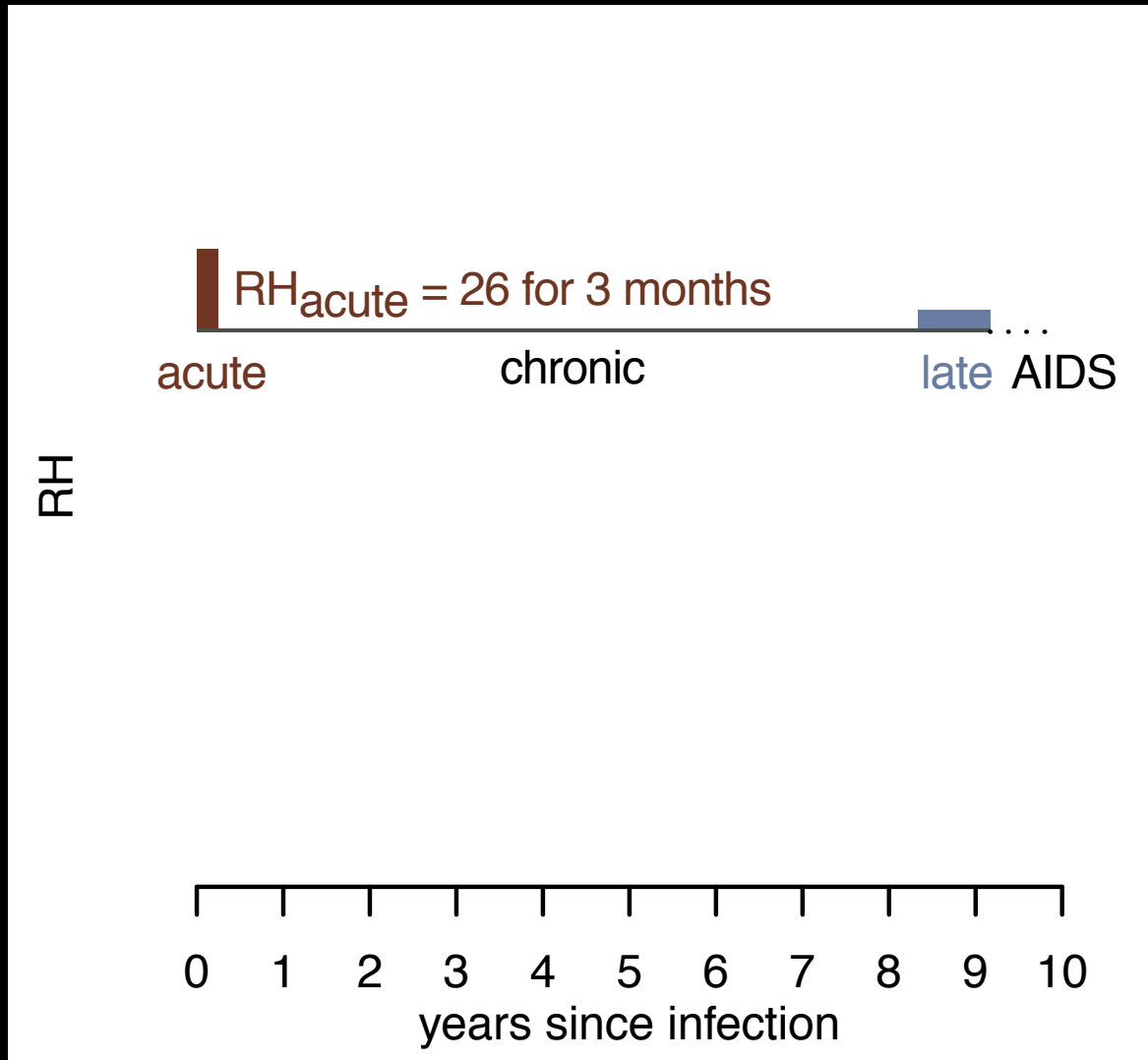


Refit the same model using Bayesian MCMC

# Collinearity in Fitted Parameters



Data are consistent with both
- shorter, highly infectious
- longer, less infectious acute phases

95% contour

95% CI

95% CI

$d_{acute}$ (months)

$RH_{acute}$

Refit the same model using Bayesian MCMC

# Collinearity in Fitted Parameters



RH$_{acute}$ = 26 for 3 months

acute    chronic    late  AIDS

RH

0  1  2  3  4  5  6  7  8  9  10
years since infection

What is actually Identifiable?

Excess Hazard-Months due to acute phase

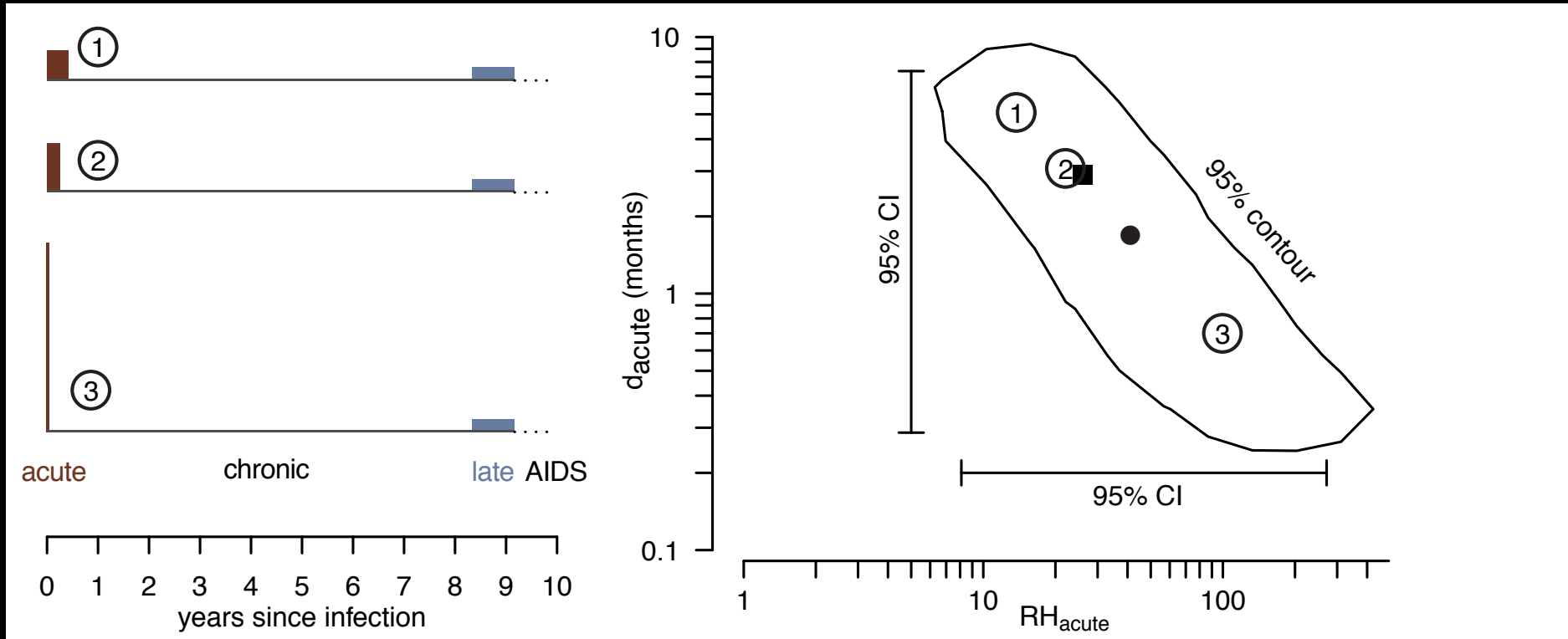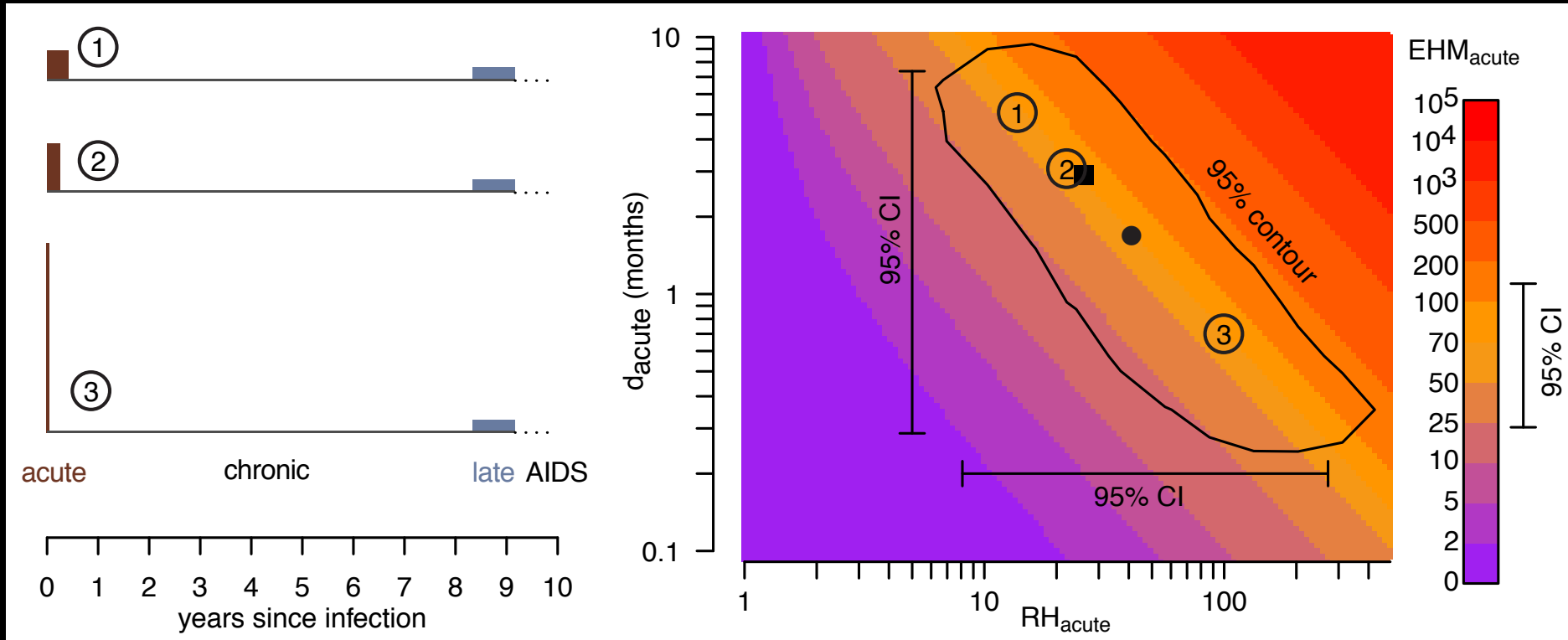EHM$_{acute}$ = (RH$_{acute}$-1)d$_{acute}$

EHM$_{acute}$ = 25*3 = 75

EHM$_{acute}$ = 15*5 = 75

EHM$_{acute}$ = 100*3/4 = 75

# Excess Hazard Months (EHM$_{acute}$)

# Excess Hazard Months (EHM$_{acute}$)



RH$_{acute}$ and d$_{acute}$ are not identifiable from 10-month interval cohorts

We should focus on EHM$_{acute}$

# Formally vs Informally Fitting

- Recently, fitting models to data expected

- Unnecessary for demonstration of qualitative dynamics

- Necessary for
  - parameter estimation
  - inference
  - formal model comparison
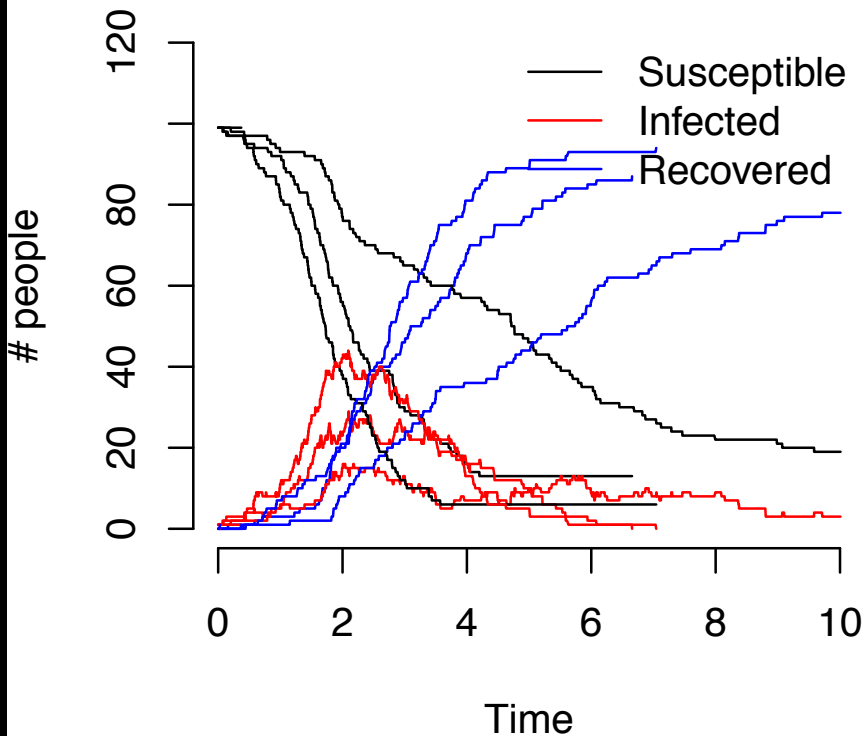
# Learning More: Methods for Fitting

- Least Squares

- Frequentist Maximum Likelihood Fitting

- Bayesian Posterior Estimation (usually MCMC)
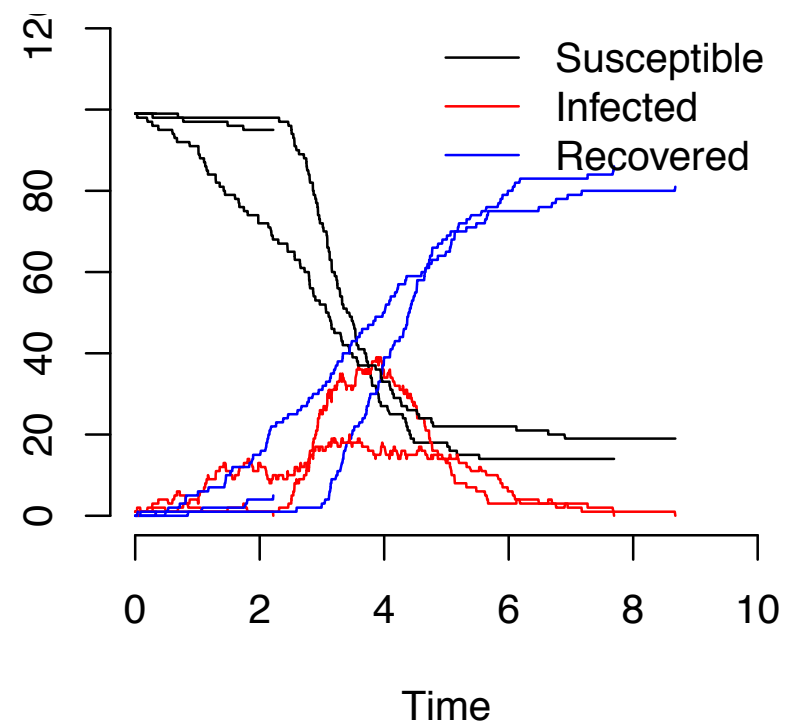
# Simulating to test methods

- Create model

- Simulate data

- Can you estimate the inputted parameters for the simulation by fitting?
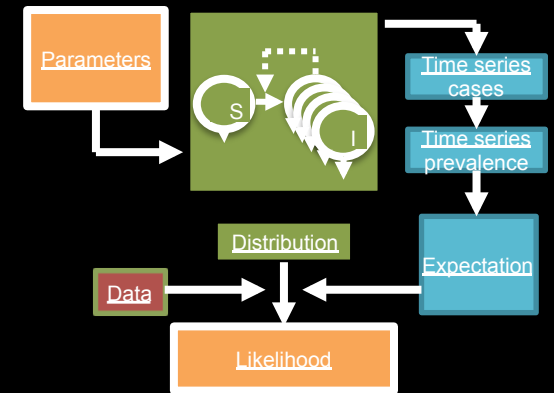
# Simulating to test methods

# Outline

1. Recap: Classical and Mechanistic Epidemiology

2. Why fit models to data?

3. Review of Linear Regression

4. Maximum Likelihood and Fitting Simple Models

5. Fitting Dynamic Models to Data

6. Summary

# Summary

- Why we fit
  parameter estimation
  inference
  formal model comparison



- How we fit
  Create a probabilistic framework that links
  our model to data—ie, write a likelihood

- What to consider when fitting
  Assumptions              Goodness of fit
  Overfitting              Identifiability

Title: Models & Data: Introduction to Model Fitting

Attribution:

Bellan SE (2012) Introduction to Model Fitting. *International Clinics on Infectious Disease Data and Dynamics*. http://ici3d.org

For further information or slides in Microsoft Powerpoint please contact Steven Bellan (steve.bellan@uga.edu).