

Model Selection and Comparison



Cara Brook, Jessica Metcalf, and Christian Ranaivoson

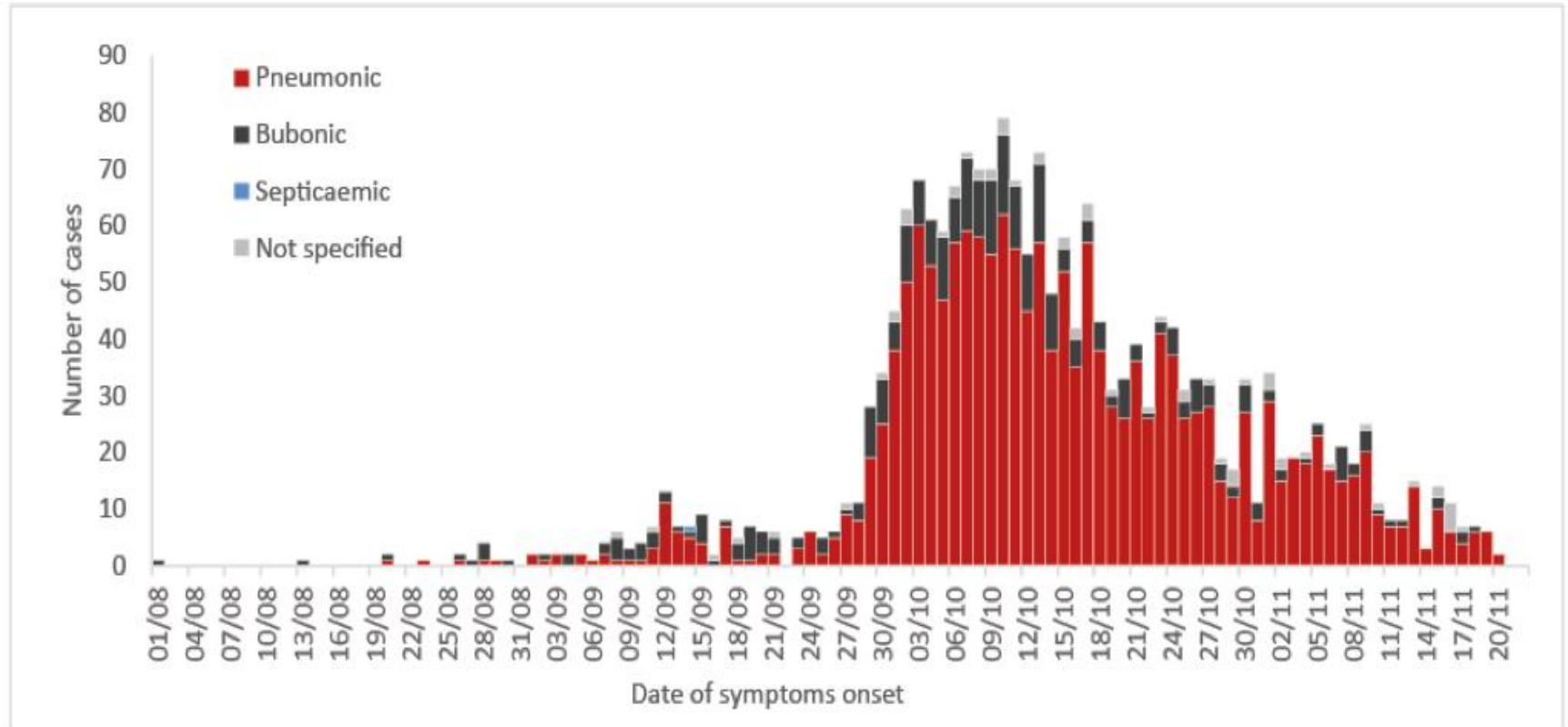
University of California, Berkeley, USA

University of Princeton

University of Antananarivo, Madagascar

E2M2 2022 Ranomafana, Madagascar

Which model is best?



There are many statistical methods used to 'fit' models to data and there are many scenarios from which mechanical model can be built.

There are many statistical methods used to 'fit' models to data and there are many scenarios from which mechanical model can be built.

The method best suited for your work will depend on your model and your data.

There are many statistical methods used to 'fit' models to data and there are many possible scenarios from which mechanical model can be built.

The method best suited for your work will depend on your model and your data.

What are some measures of model fit used in E2M2 so far?

R squared

R-carré

Least squares

(Moindres carrés)

Log likelihood

Maximum de vraisemblance

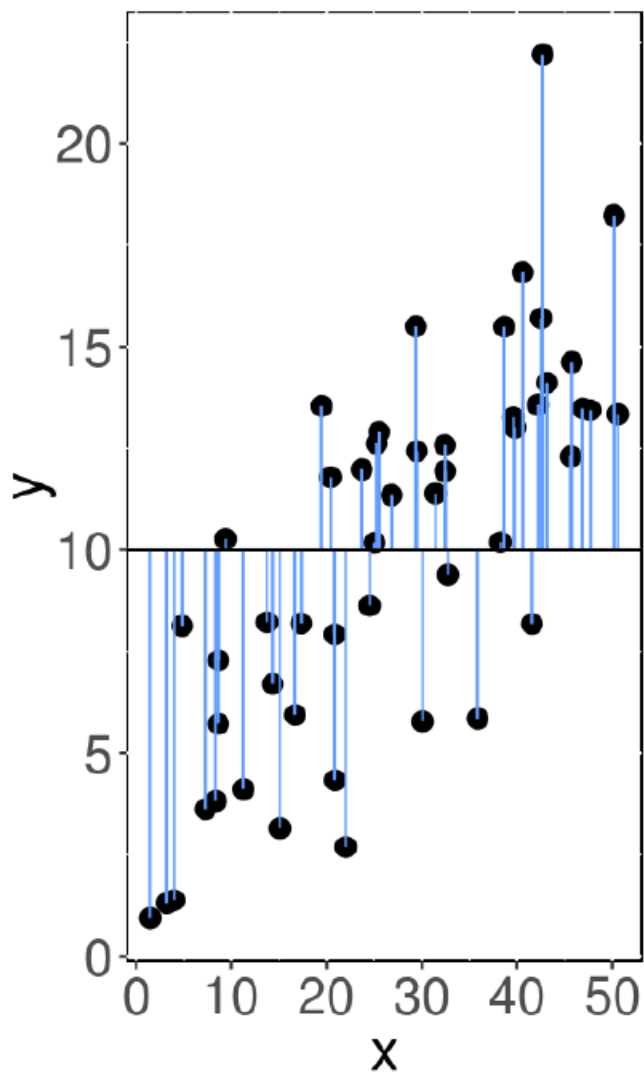
AIC

*(uses least squares or log-likelihood but penalizes by
number of fitted parameters)*

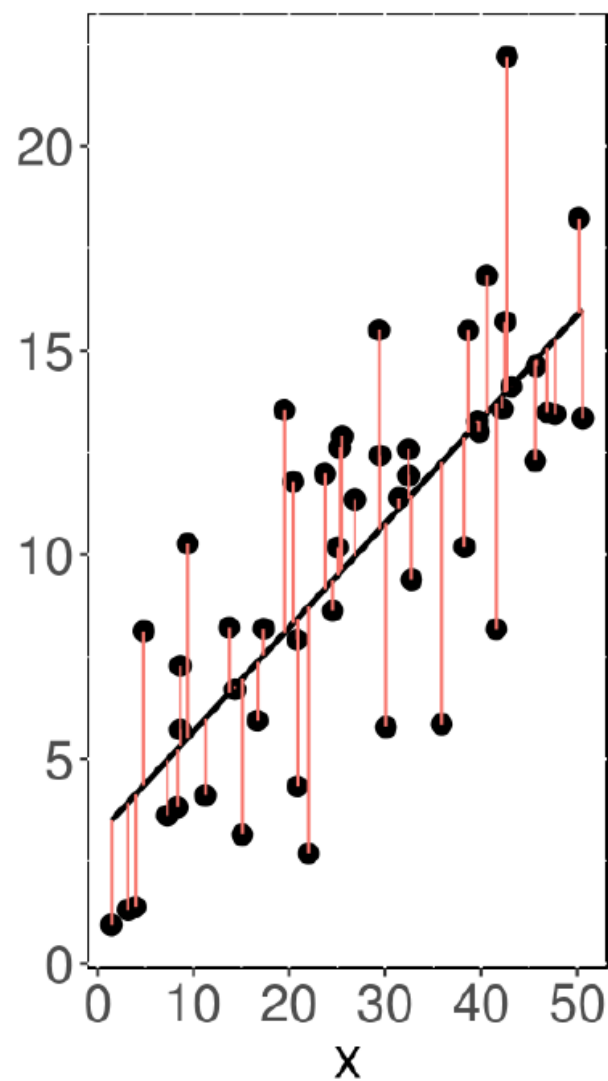


Definition r^2

SS total



SS error

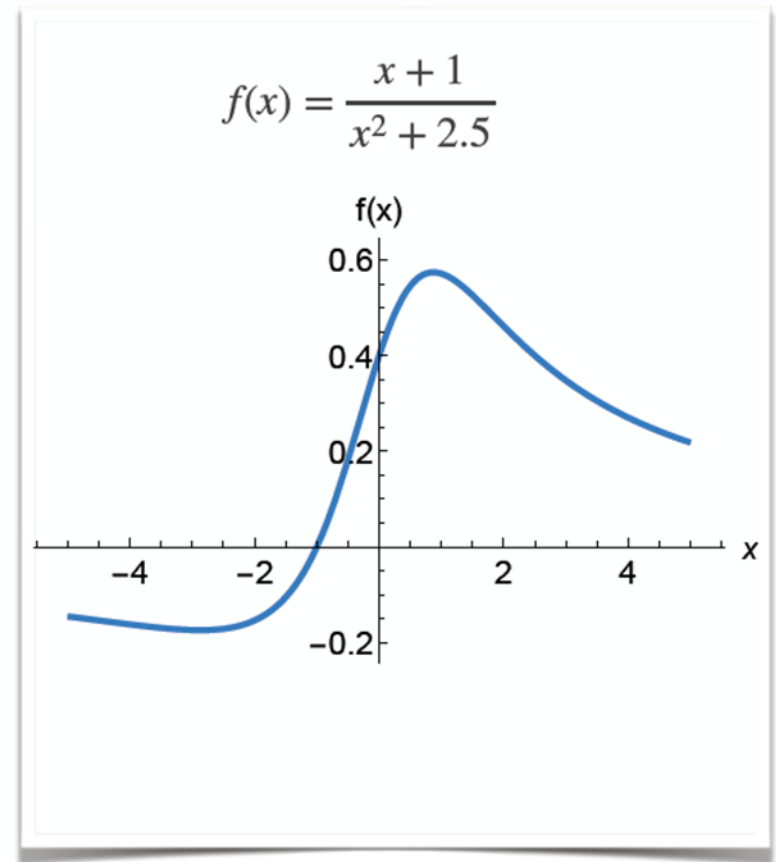


$$R^2 = 1 - \frac{SSE_p}{SST}$$

Optimization/maximization

Function properties

?

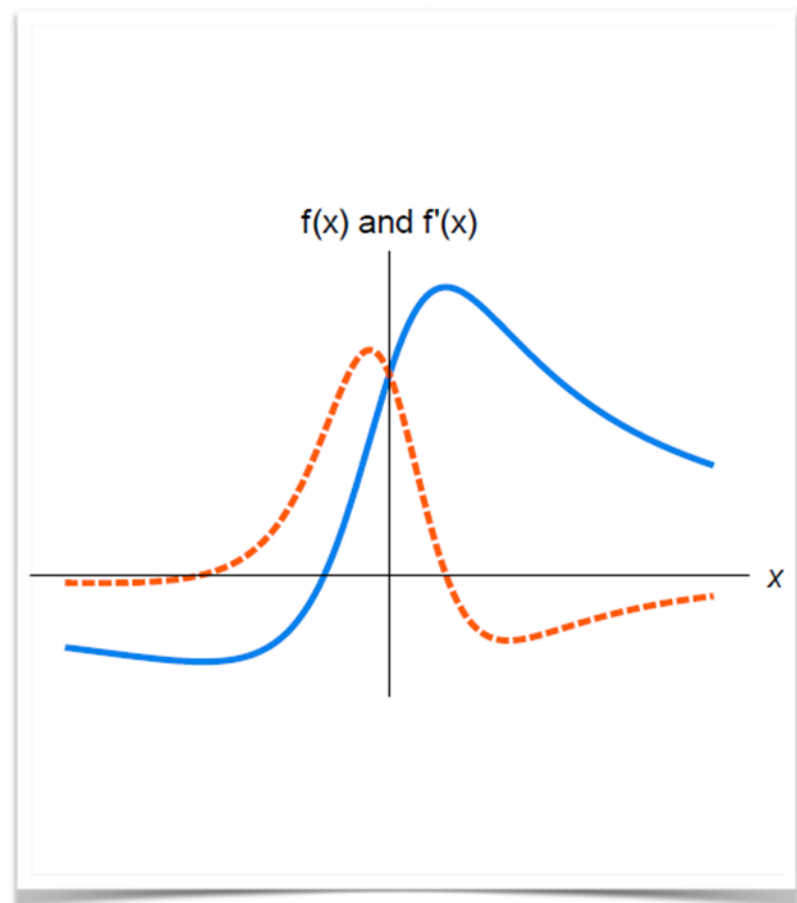


From Tanjona Ramiadantsoa

Optimization/maximization

A function and its derivative

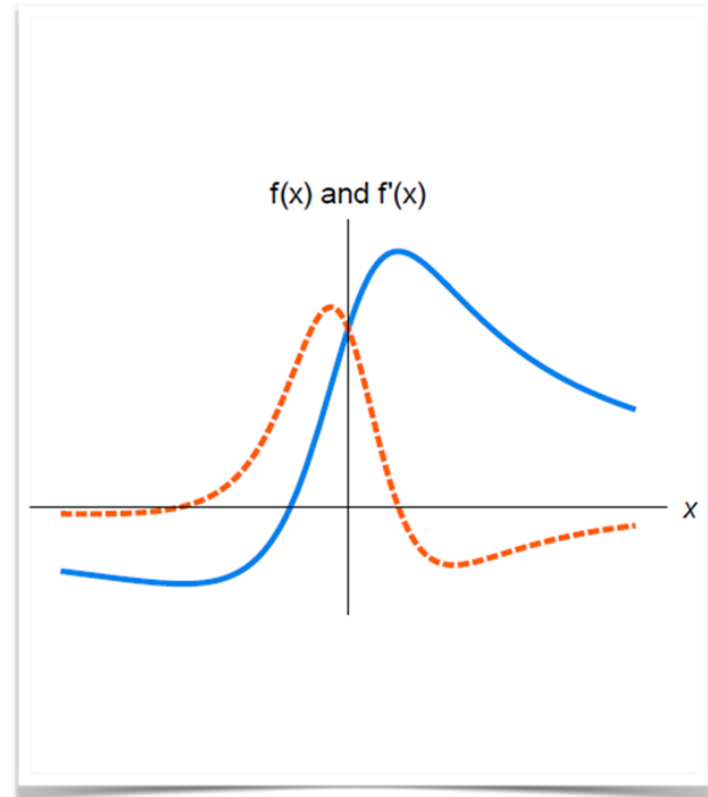
- ❖ What happens when the derivative is:
 - ❖ negative?
 - ❖ positive?
 - ❖ zero?
 - ❖ reaching a maximum (finite) value?



Optimization/maximization

A function and its derivative

- ❖ What happen when the derivative is:
 - ❖ negative?
 - ❖ positive?
 - ❖ zero?
 - ❖ reaching a maximum (finite) value?

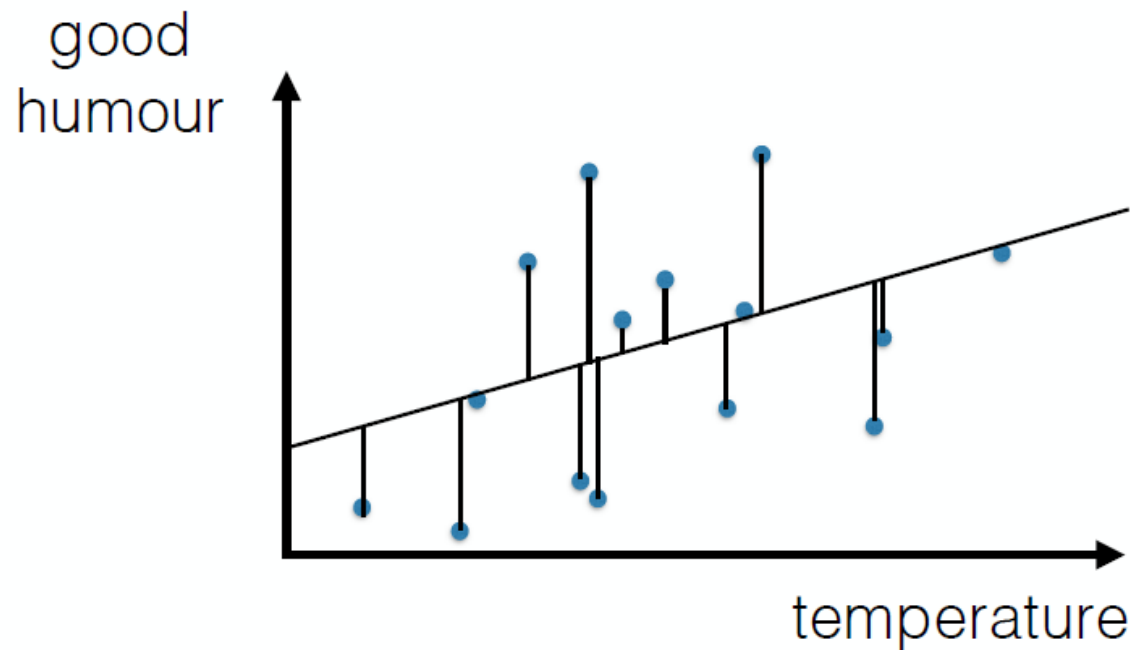


From Tanjona Ramiadantsoa

The R function 'optim' can be used to minimize these measures of model difference from the data.

Least squares

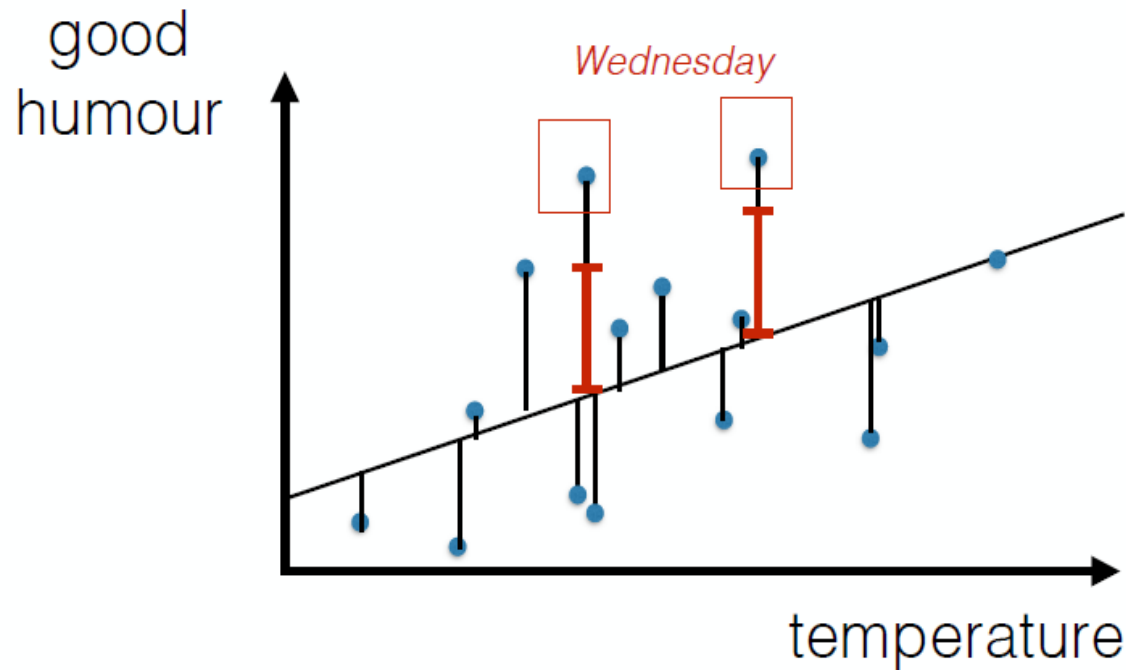
Adding covariates and R^2



$$\text{humour} = b_0 + b_1 \text{temperature} + \text{Error}$$

Least squares

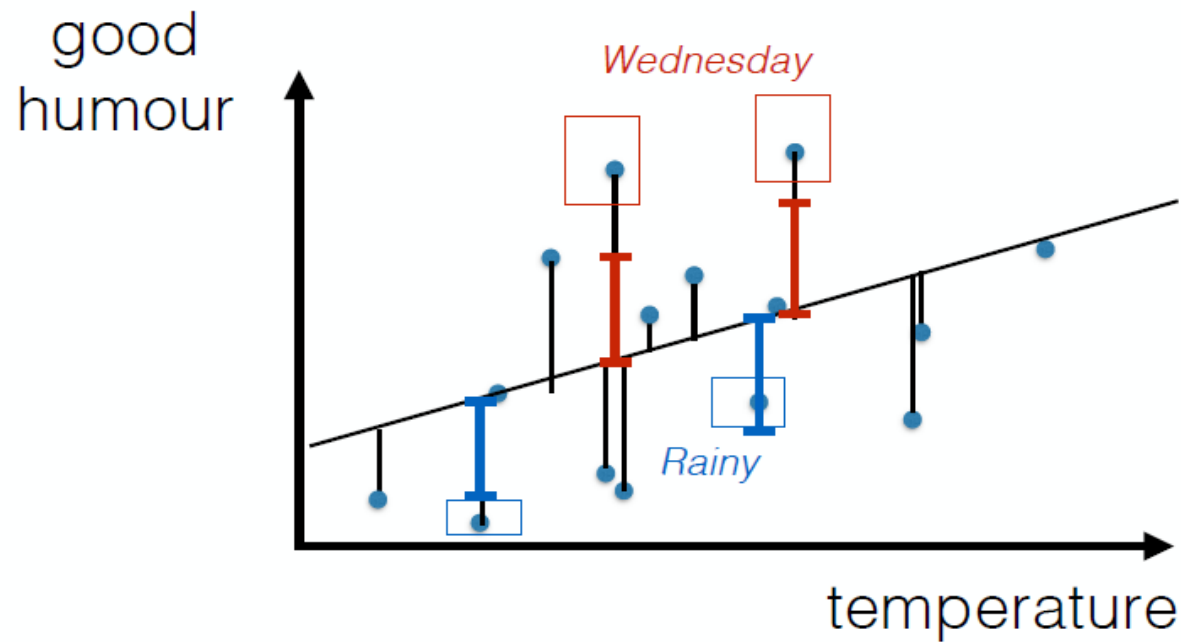
Adding covariates and R^2



$$\text{humour} = b_0 + b_1\text{temperature} + b_2\text{Wednesday} + \text{Error}$$

Least squares

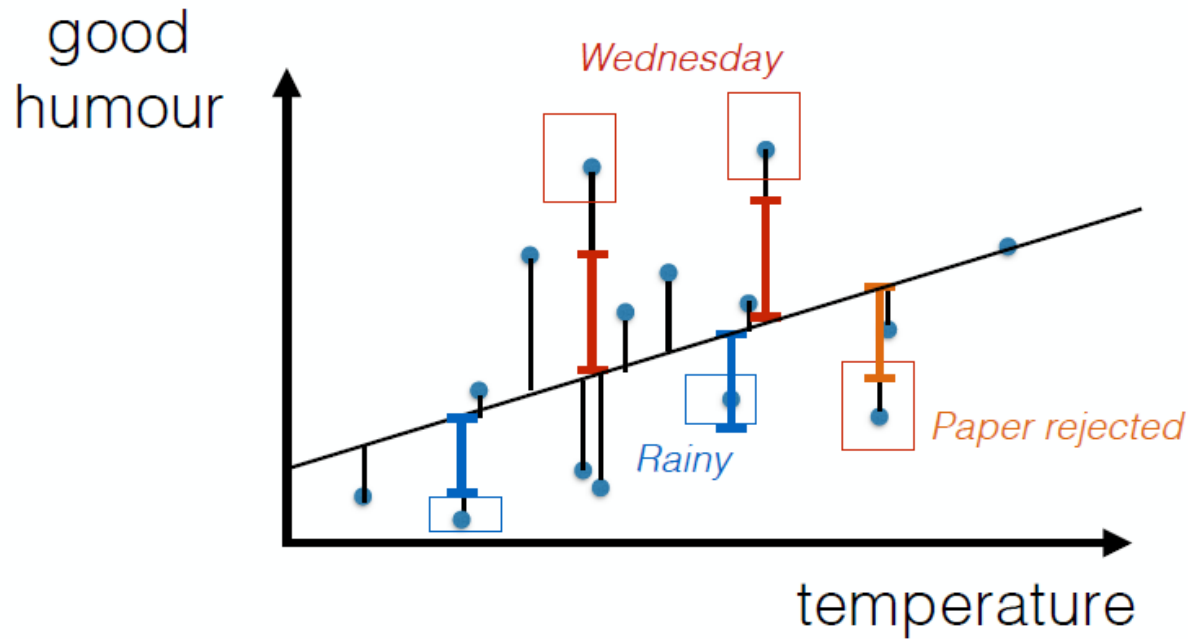
Adding covariates and R^2



$$\text{humour} = b_0 + b_1\text{temperature} + b_2\text{Wednesday} + b_3\text{rain} + \text{Error}$$

Least squares

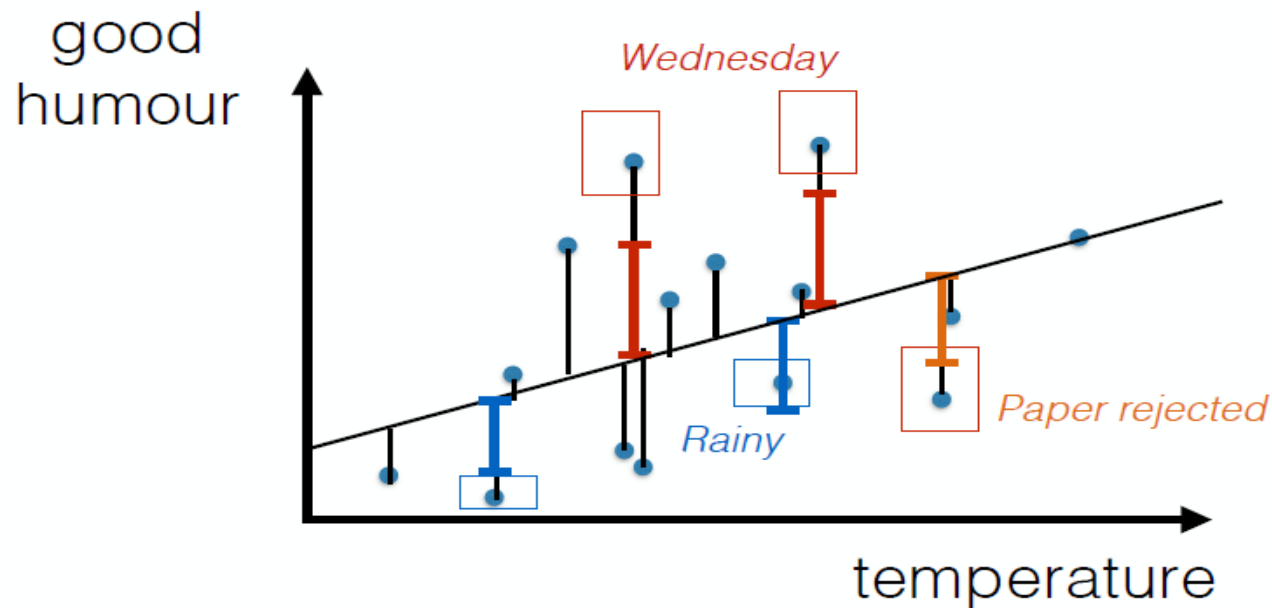
Adding covariates and R^2



$$\text{humour} = b_0 + b_1\text{temperature} + b_2\text{Wednesday} + b_3\text{rain} + b_4\text{rejection} + \text{Error}$$

Least squares

Adding covariates and R^2



Adding covariates almost always increases the R^2 - so a key question is when to stop.

What to choose?



Least square AIC

$$AIC = N * \ln\left(\frac{SS_e}{N}\right) + 2K$$

***N**: Number of observations*


***SS_e**: Sum square of errors*

***K**: Number of parameters*

The smaller the AIC the better

Least square AIC

More parameter is not always good


$$\text{AIC} = N * \ln\left(\frac{SS_e}{N}\right) + 2K$$

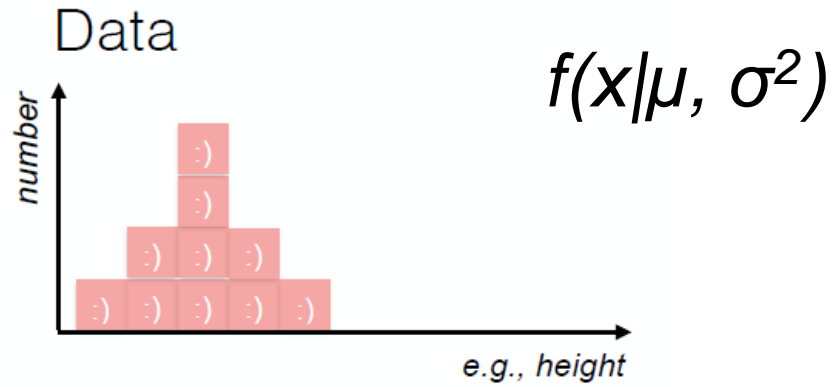
***N**: Number of observations*

***SS_e**: Sum square of errors*

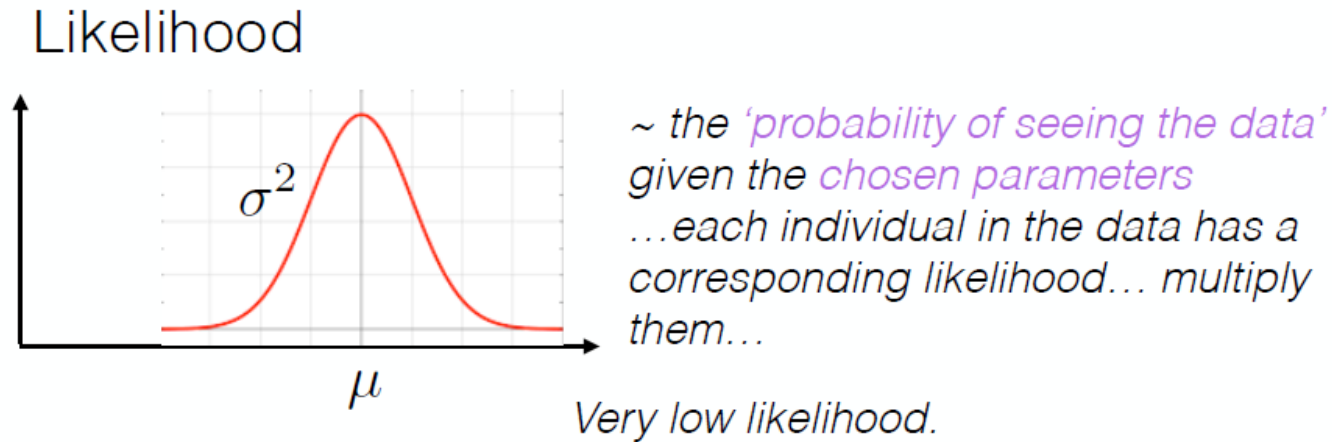
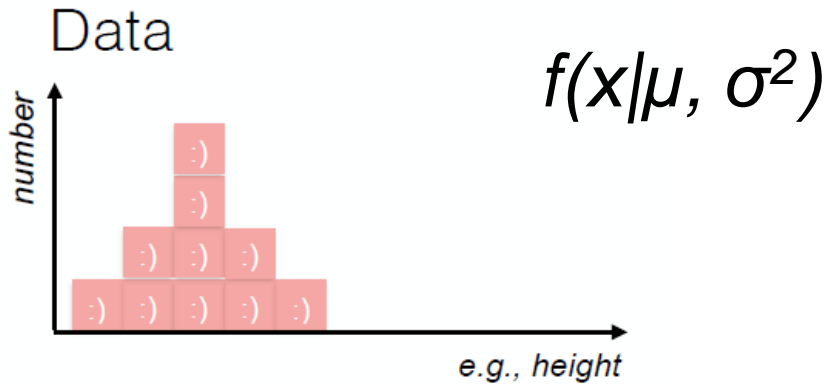
***K**: Number of parameters*

The smaller the AIC the better

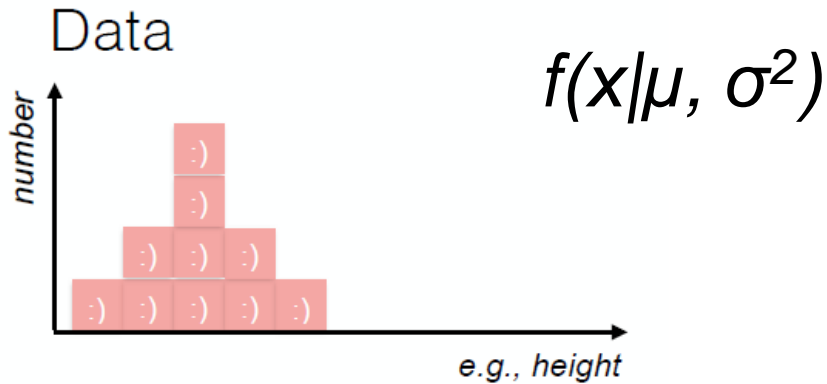
Maximum likelihood



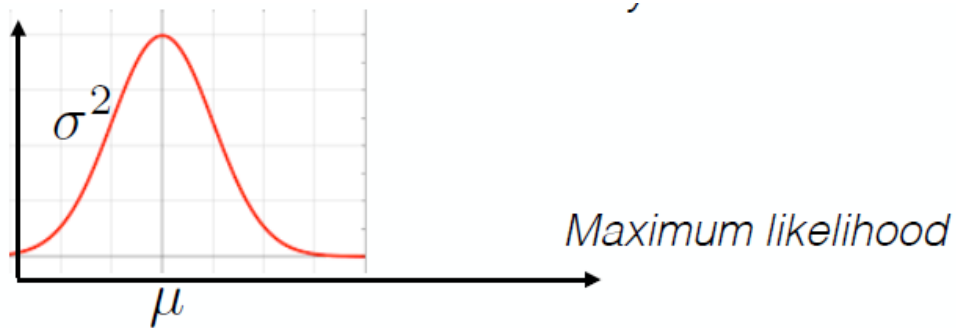
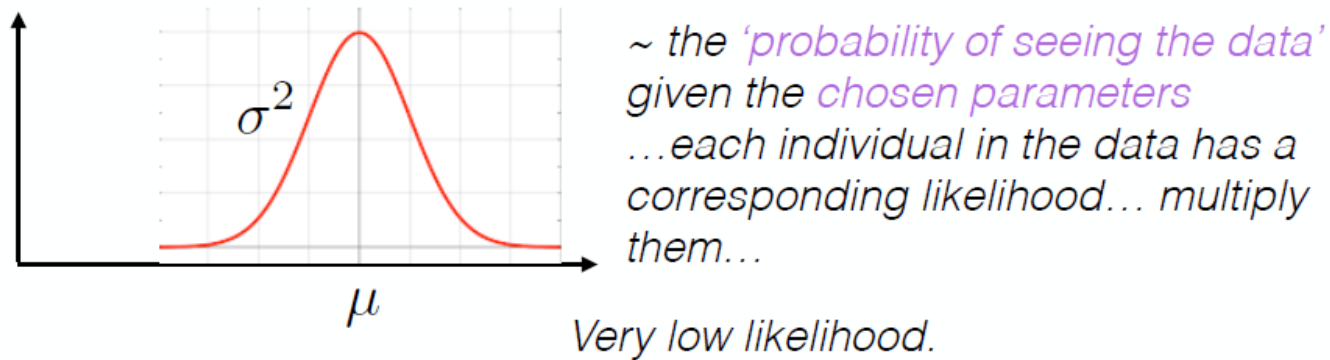
Maximum likelihood



Maximum likelihood



Likelihood



An example of model selection: *Bartonella* spp. in Madagascar rats

Epidemics 20 (2017) 56–66



Contents lists available at ScienceDirect

Epidemics

journal homepage: www.elsevier.com/locate/epidemics



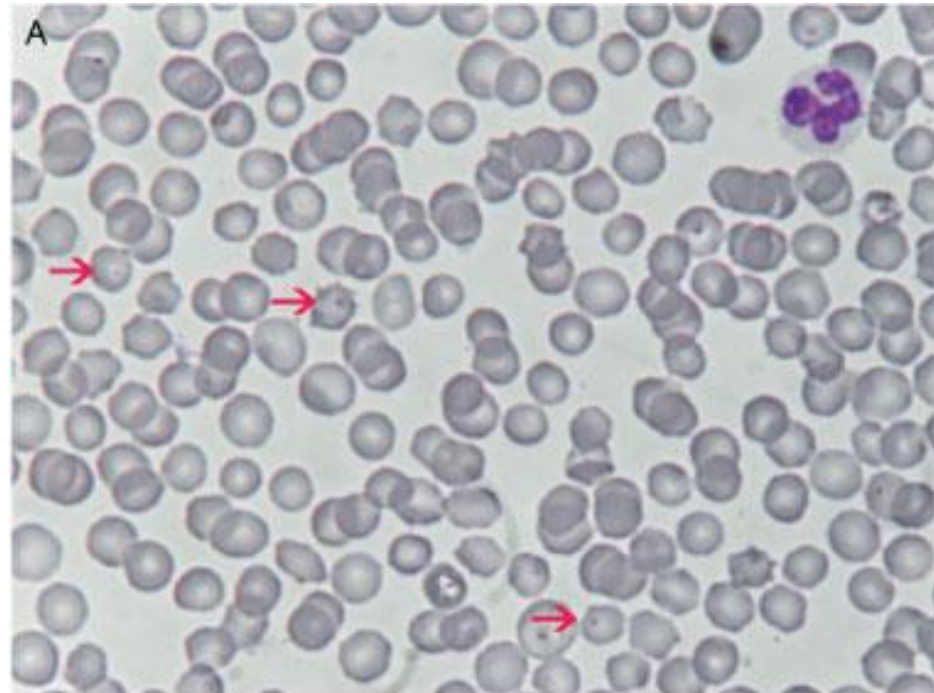
Elucidating transmission dynamics and host-parasite-vector relationships for rodent-borne *Bartonella* spp. in Madagascar



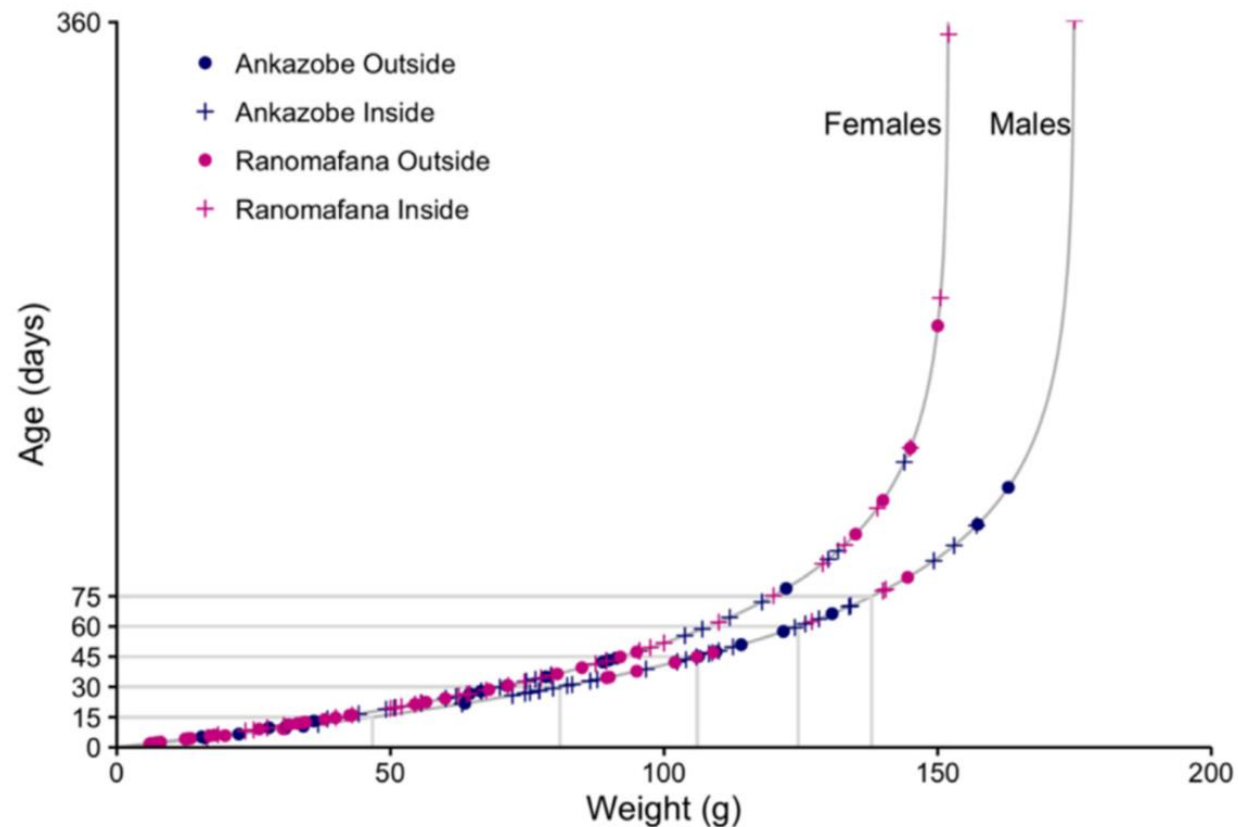
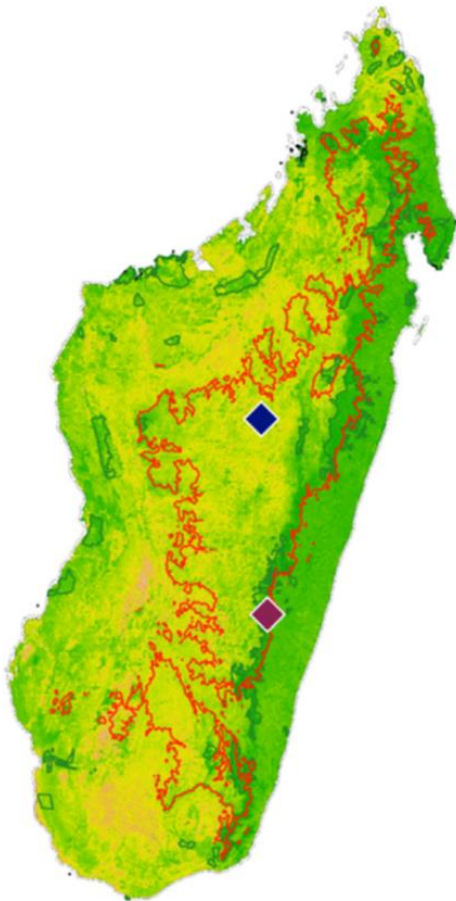
Cara E. Brook^{a,*}, Ying Bai^b, Emily O. Yu^a, Hafaliana C. Ranaivoson^{c,d}, Haewon Shin^e, Andrew P. Dobson^a, C. Jessica E. Metcalf^{a,1}, Michael Y. Kosoy^{b,1}, Katharina Dittmar^{e,1}

Bartonella spp.

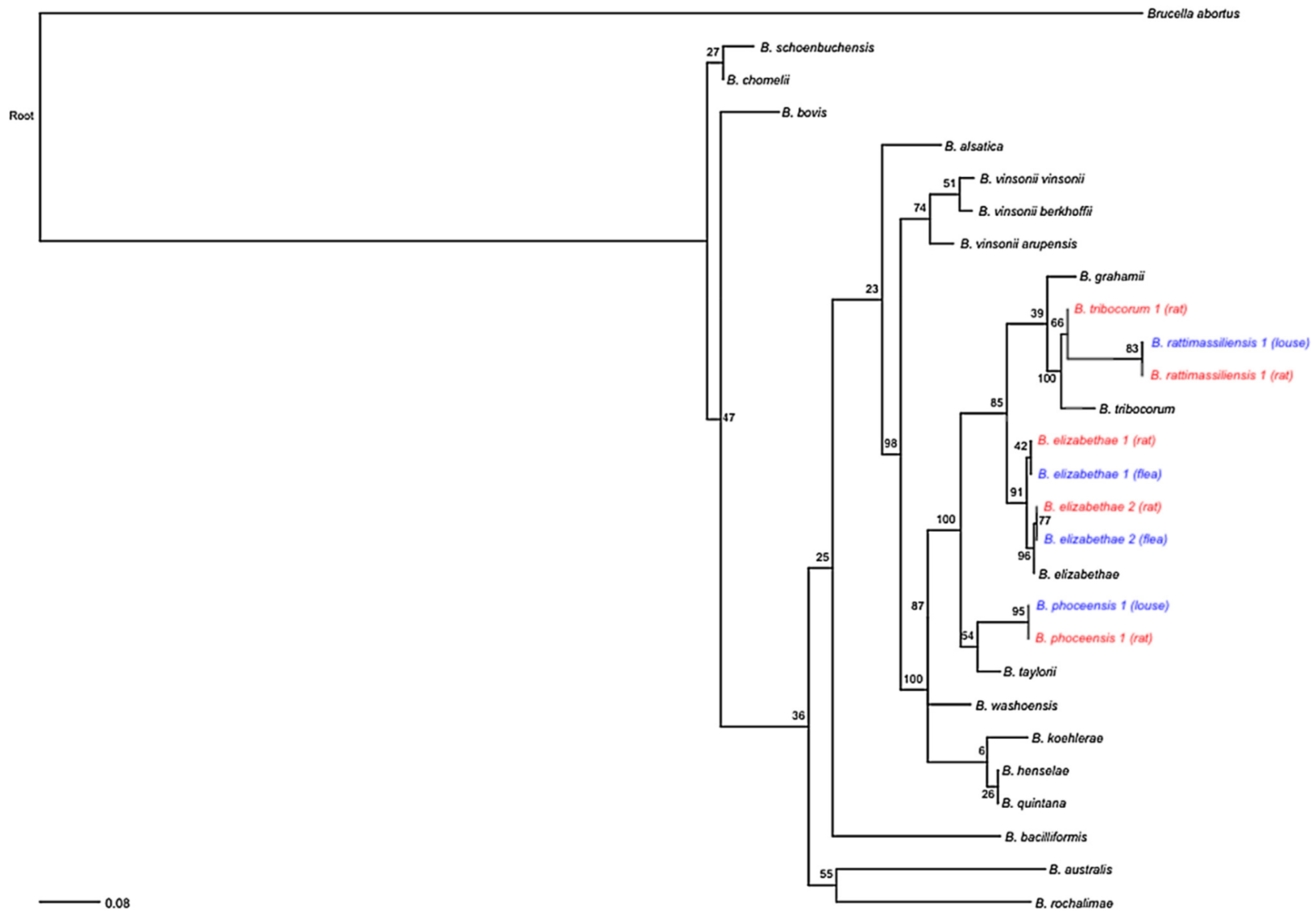
- persistent erythrocytic bacteria that are sometimes zoonotic
- vectored by ticks, fleas, sand flies, mosquitoes
- at least 8 human-infecting species
 - *Bartonella bacilliformis* = Carrion's disease
 - *Bartonella henselae* = cat scratch fever
 - *Bartonella quintana* = trench fever



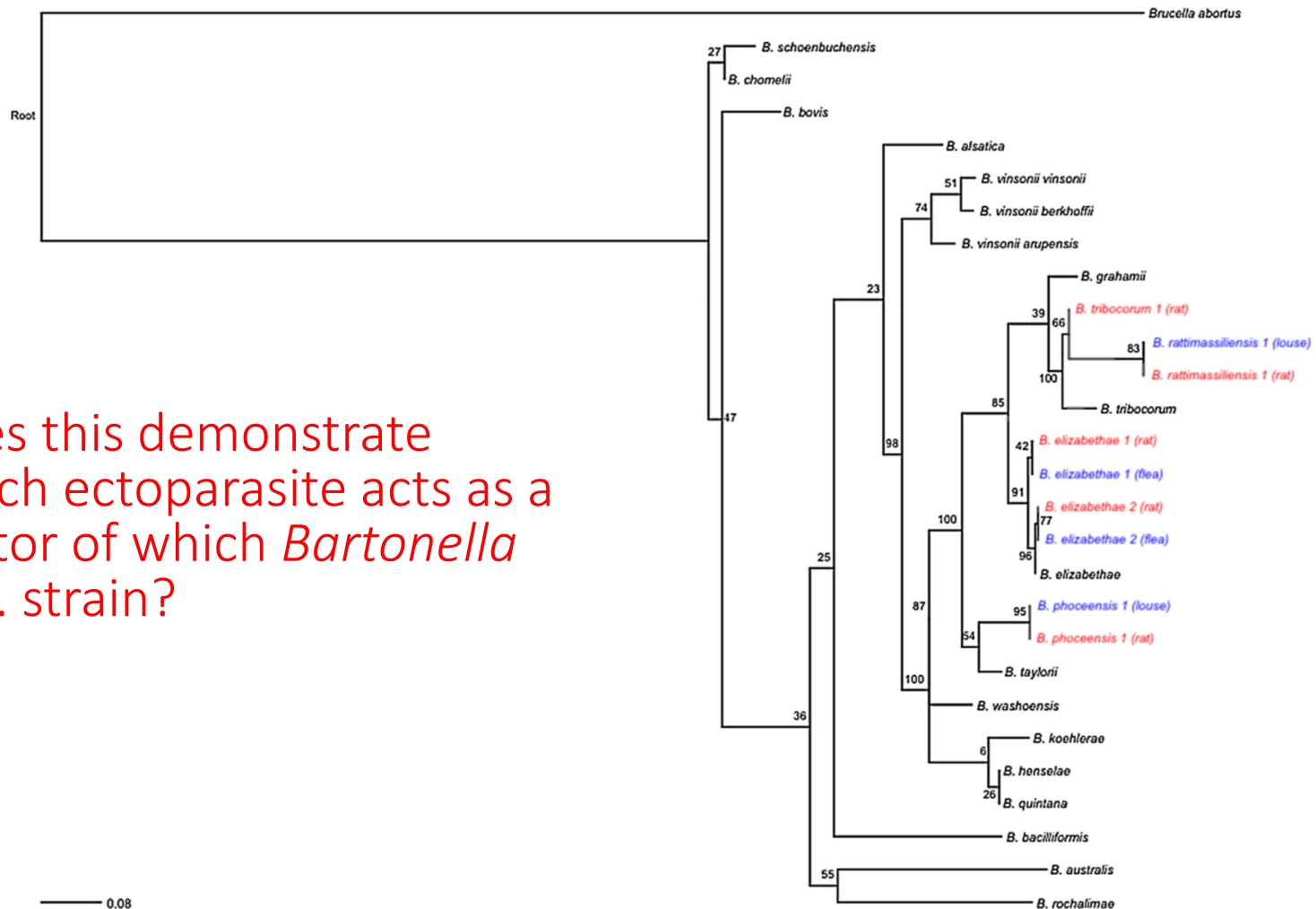
We first collected samples from rats in two sites Madagascar.



Statistically, we demonstrated an association between genotypes of *Bartonella* spp. found in rats and their ectoparasites.



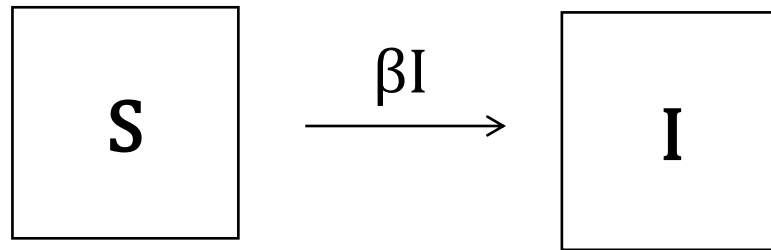
Statistically, we demonstrated an association between genotypes of *Bartonella* spp. found in rats and their ectoparasites.



Does this demonstrate
which ectoparasite acts as a
vector of which *Bartonella*
spp. strain?

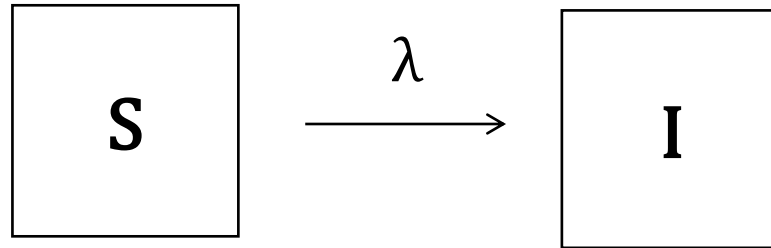
Then, we asked:
*How does the rate of becoming
infected vary with age?*

Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.



for a persistent, non-immunizing infection

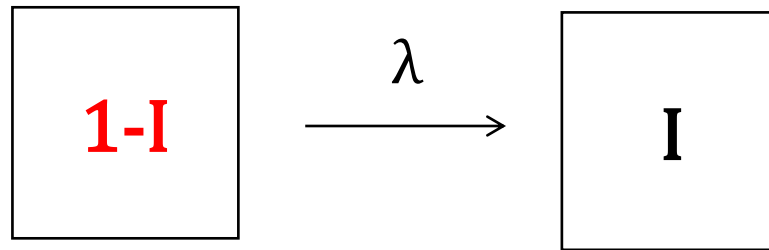
Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.



where λ , the force of infection, is the per capita rate at which susceptible hosts become infected

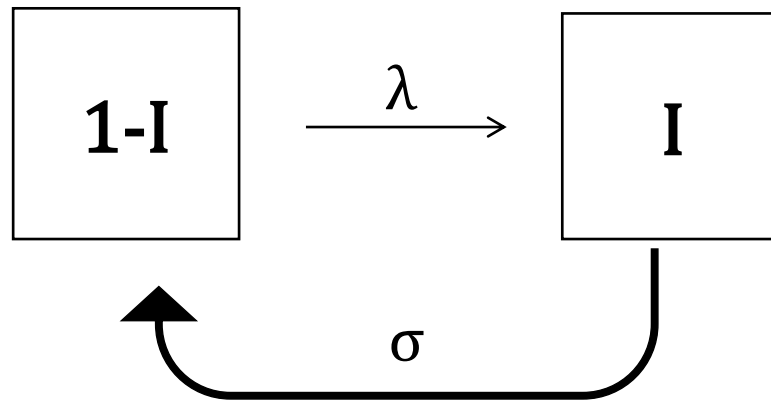
Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.

with a persistent infection,
we can assume that, if not
infected, you must be
susceptible....



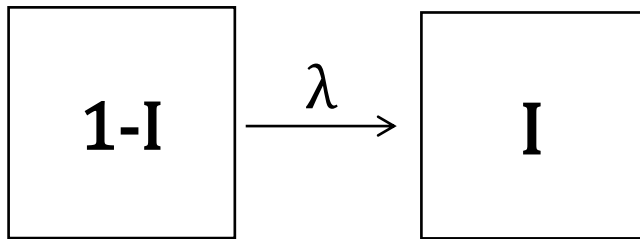
where λ , the force of infection, is the per capita rate at which susceptible hosts become infected

Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.

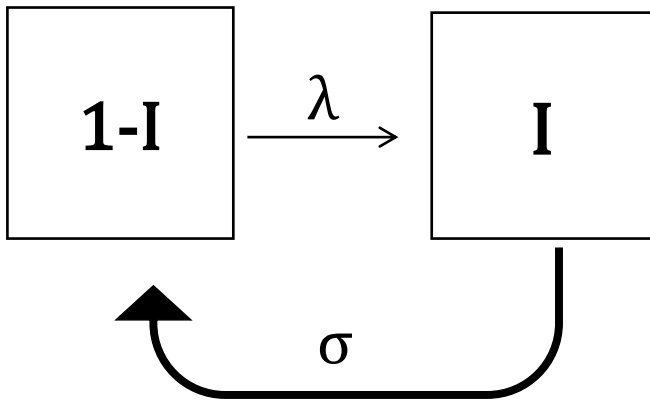


and σ is the rate of recovery from infection

Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.

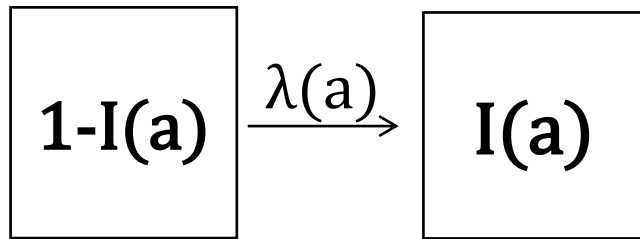


$$\frac{dI(a)}{da} = \lambda(a)(1 - I(a))$$

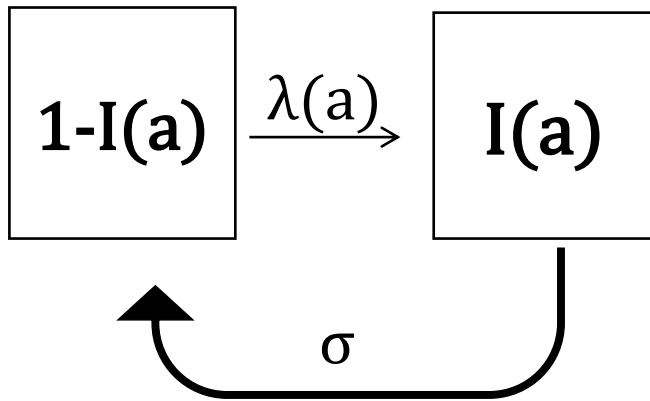


$$\frac{dI(a)}{da} = \lambda(a)(1 - I(a)) - \sigma I(a)$$

Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.

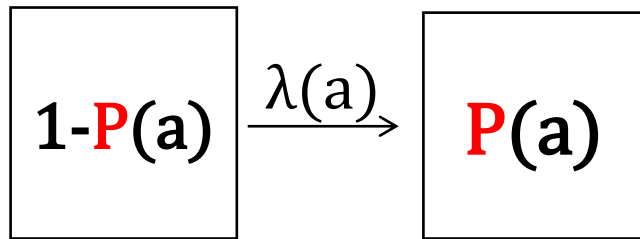


$$\frac{dI(a)}{da} = \lambda(a)(1 - I(a))$$

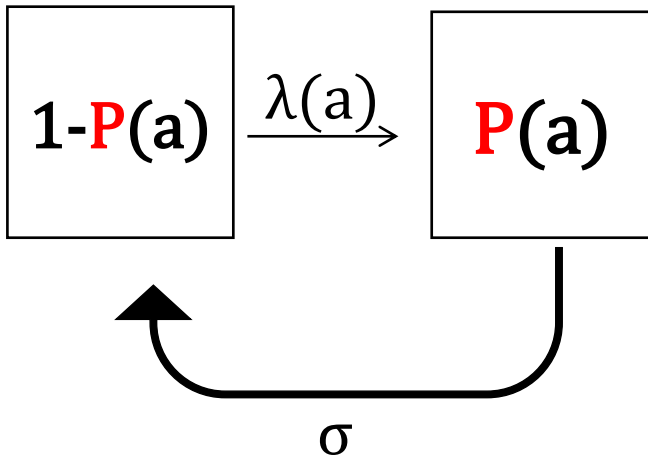


$$\frac{dI(a)}{da} = \lambda(a)(1 - I(a)) - \sigma I(a)$$

Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.



$$\frac{d P(a)}{da} = \lambda(a) (1 - P(a))$$



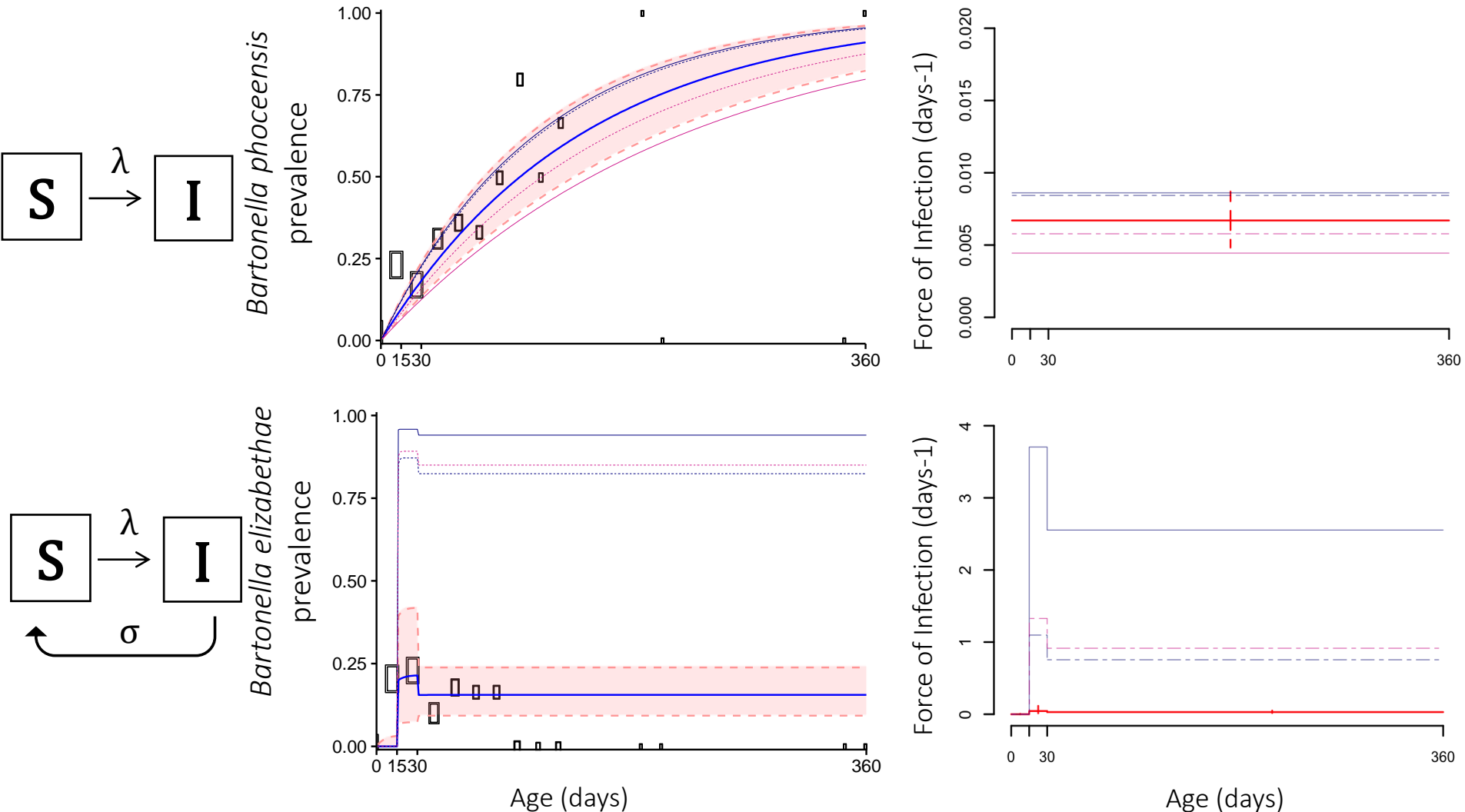
$$\frac{d P(a)}{da} = \lambda(a) (1 - P(a)) - \sigma(a) P(a)$$

Compare using **$AIC = 2K - 2\ln(L)$**

similar techniques can also be applied to age-seroprevalence data for immunizing infections

Let's see which model works best
for your data!

We found that an **SI model** offered the best fit to *B. phoceensis* data while the **SIS model** offered the best fit to the *B. elizabethae* data.



The age-structured FOI identifies age cohorts most influential in an epidemic. Juveniles showed the highest FOI.

