

# Model Selection and Comparison



Cara Brook, Jessica Metcalf, and Christian Ranaivoson

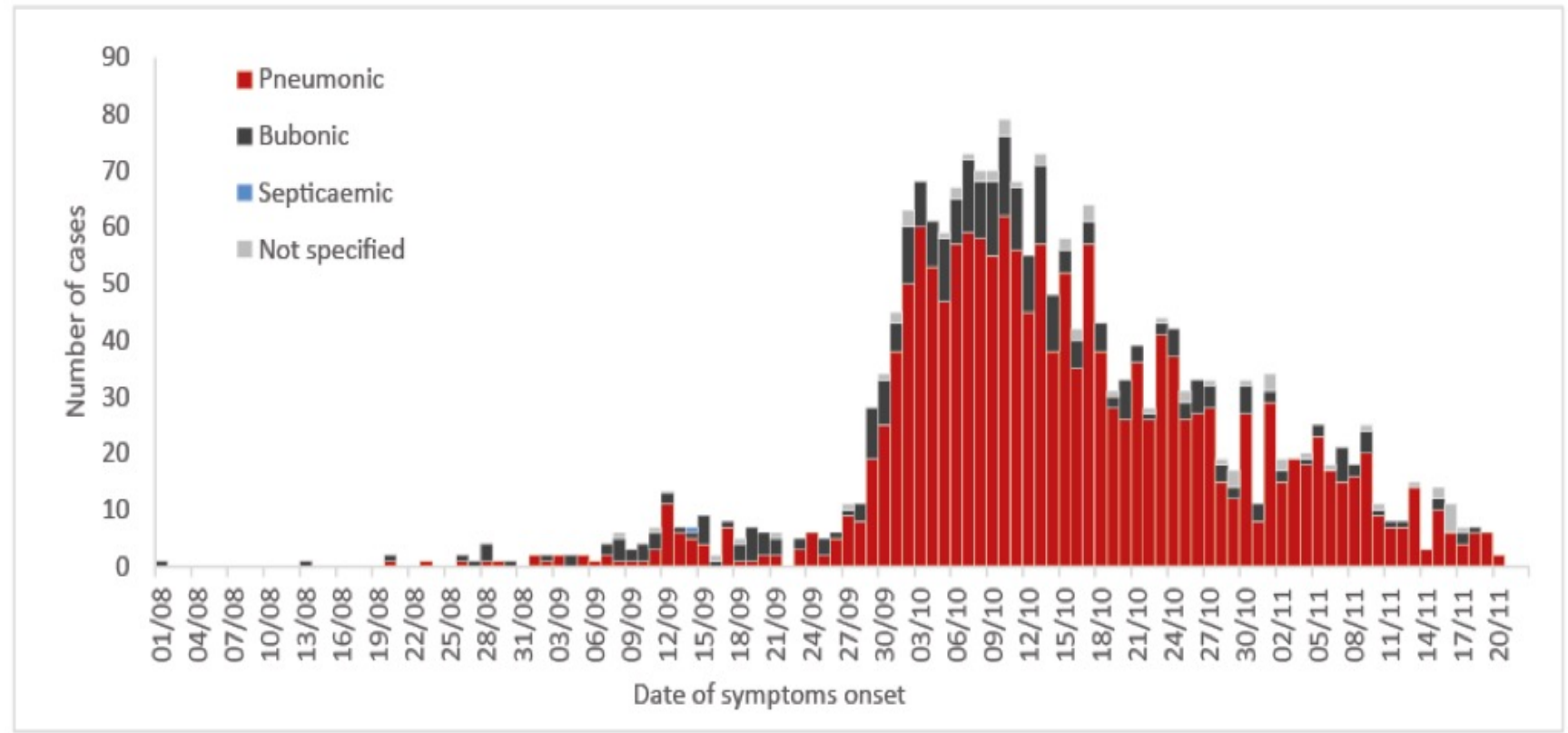
University of California, Berkeley, USA

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 possible scenarios from which mechanical model can be built.

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.

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 that you could use?

What are some measures of model fit that you could use?

R-squared

(R-carré)

Least squares

(Moindres carrés)

Maximum likelihood

(Maximum de vraisemblance)

(manakaiky indrindra ny tena izy)

AIC

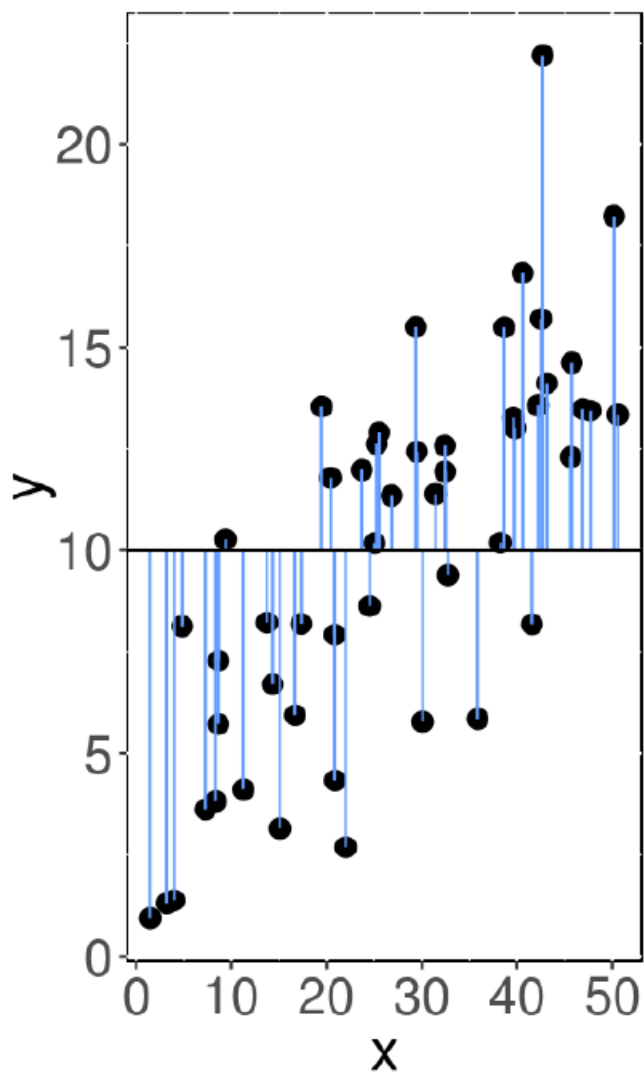
*(uses least squares or log-likelihood but penalizes by number of fitted parameters)*

Hirotsugu Akaike

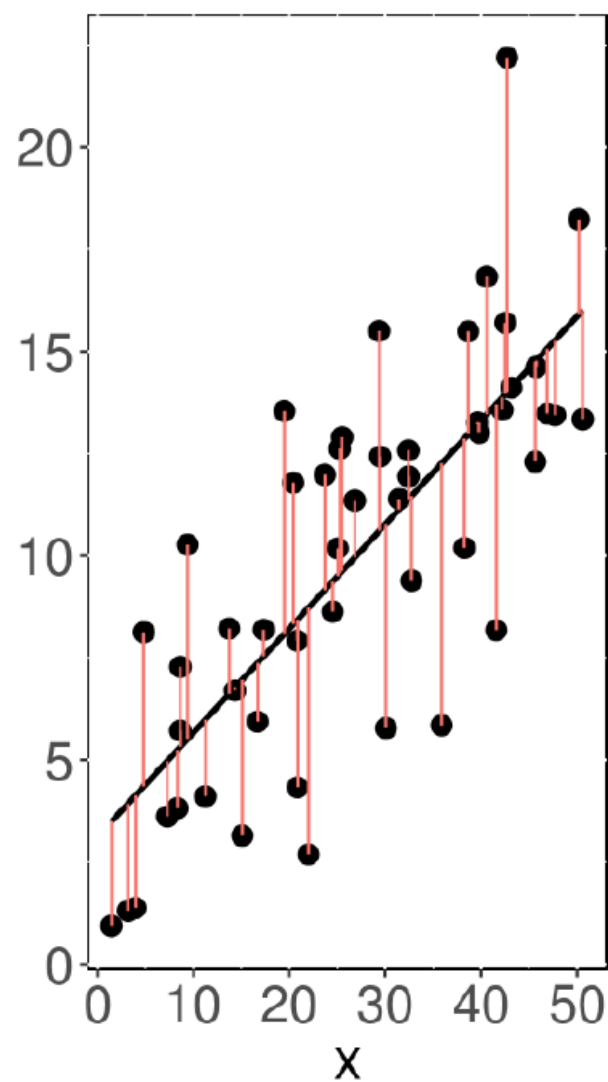


## Definition $r^2$

SS total



SS error



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

# Likelihood      Vraisemblance

$$L(\theta) = \prod_{i=1}^n f(x_i|\theta)$$

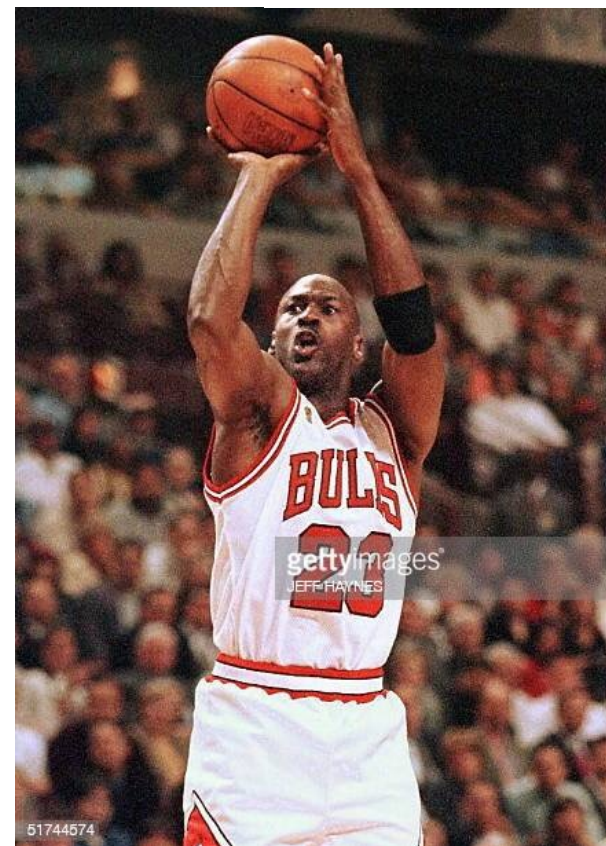
$$l(\theta|x) = \log L(\theta|x)$$

Examples:

R function : **dbinom(x, size, prob, log=T)**

➡ R function : **dbinom(8, 10, 0.835) = ?**

SUMMARY	FG3%	FT%
Career	32.7	83.5





# Likelihood      Vraisemblance

$$L(\theta) = \prod_{i=1}^n f(x_i|\theta)$$

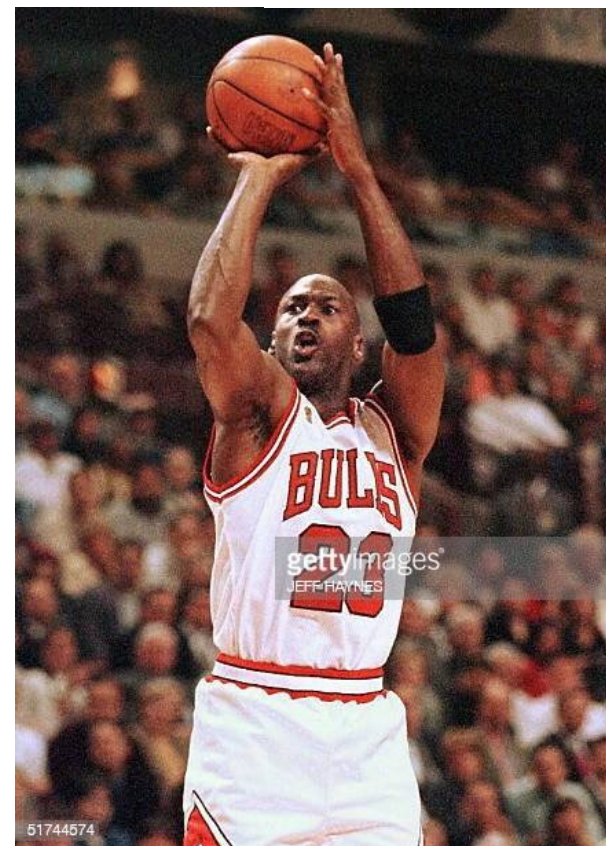
$$l(\theta|x) = \log L(\theta|x)$$

Examples:

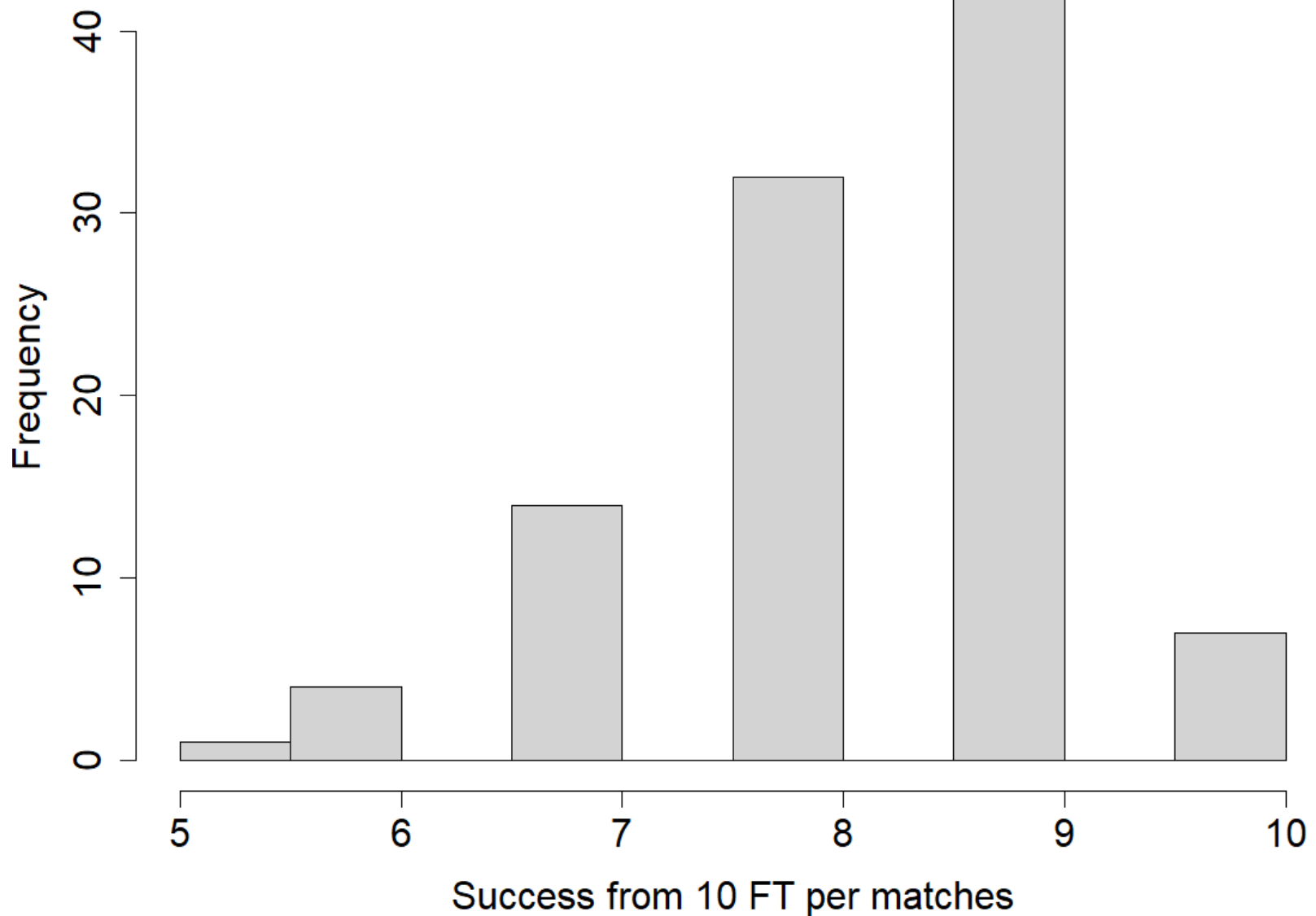
R function : **`dbinom(x, size, prob, log=T)`**

➡ R function : **`dbinom(8, 10, 0.835) = 0.289`**

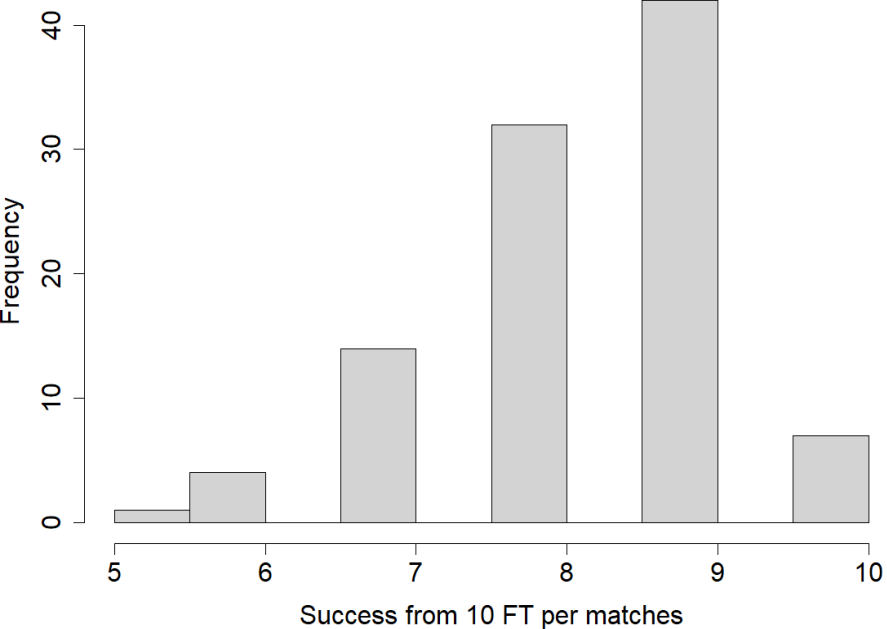
SUMMARY	FG3%	FT%
Career	32.7	83.5



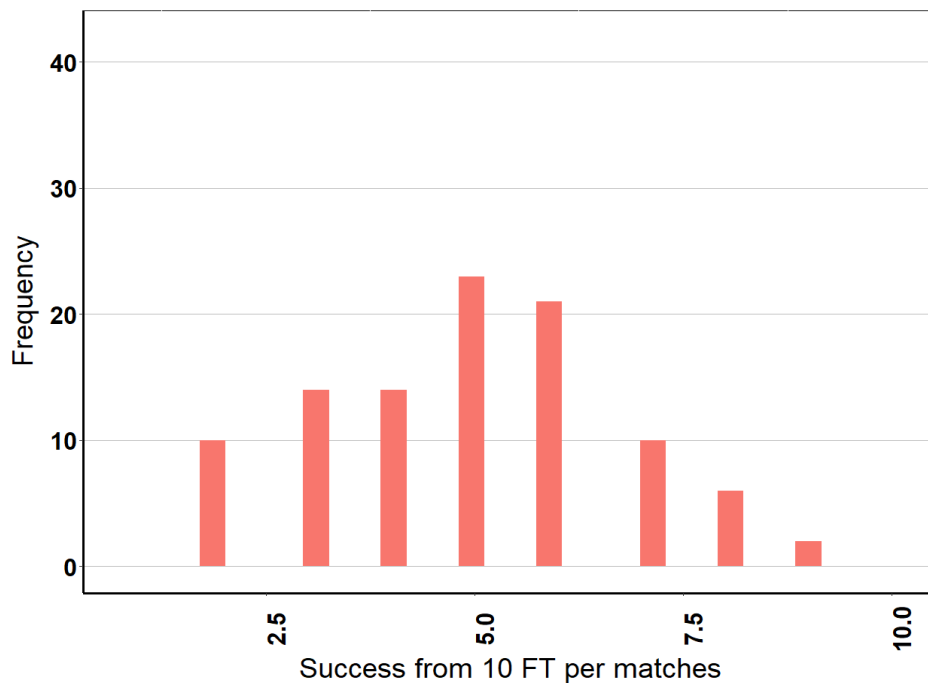
# Michael Jordan's FT success in 100 match career



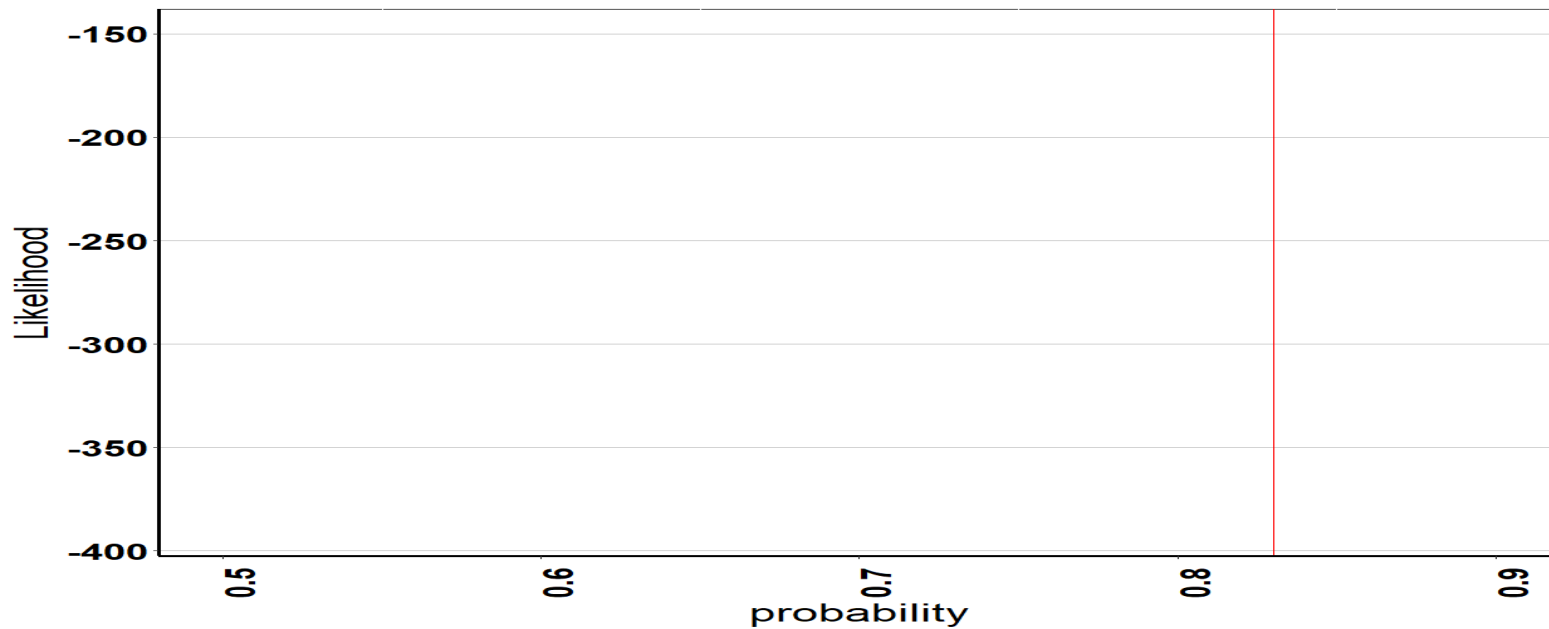
Michael Jordan's FT success in 100 match career



Michael Jordan's FT success in 100 match career



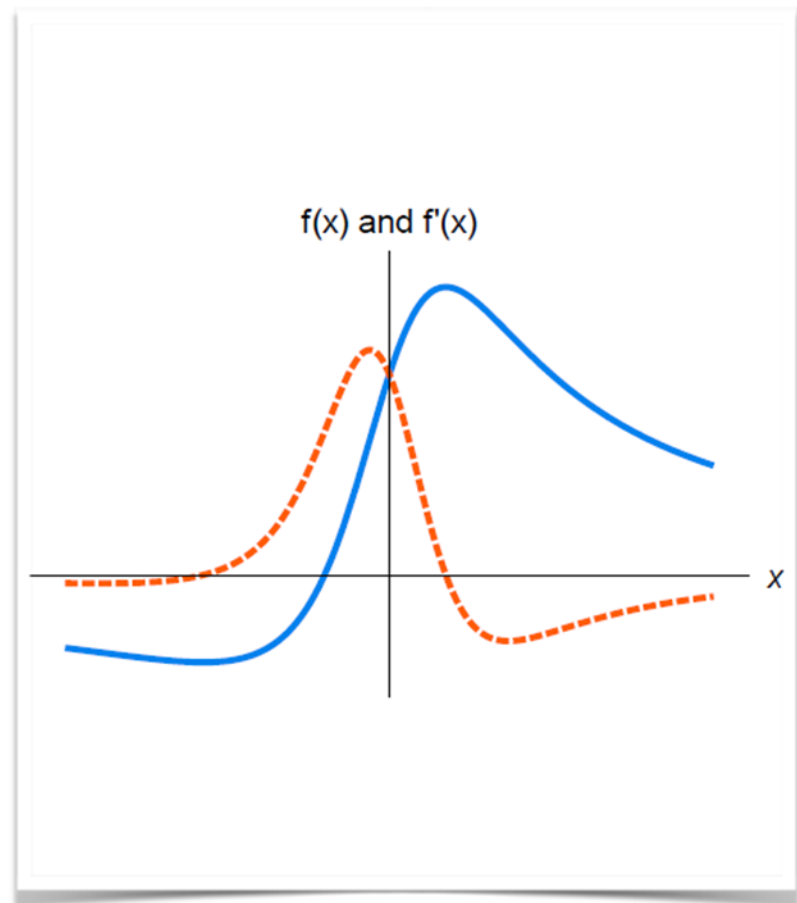
Michael Jordan's FT success in 100 match career



# Optimization/maximization

## A function and its derivative

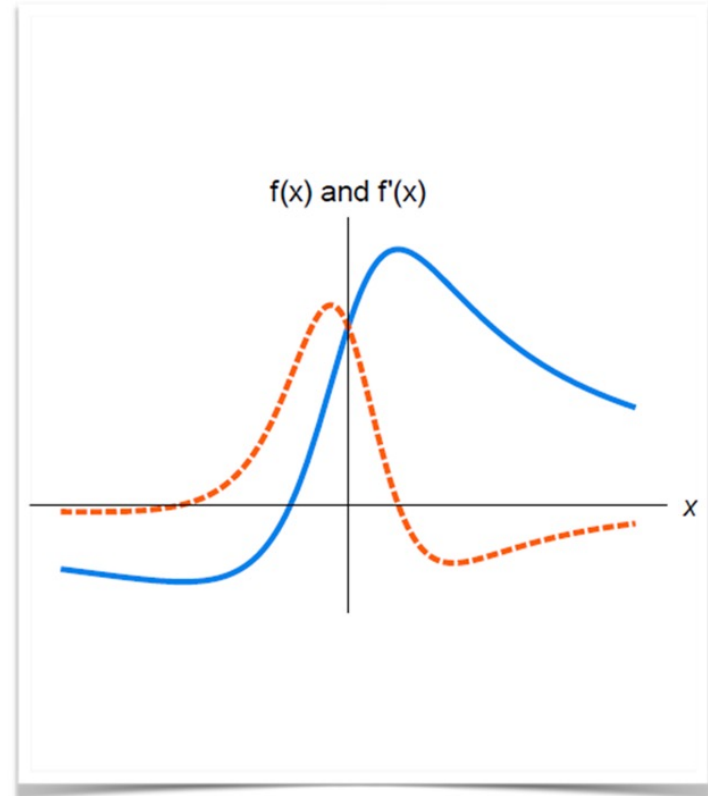
- ❖ What happens when the derivative is:
  - ❖ negative?
  - ❖ positive?
  - ❖ zero?
  - ❖ reaching a maximum (finite) value?



# Optimization/maximization

## A function and its derivative

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

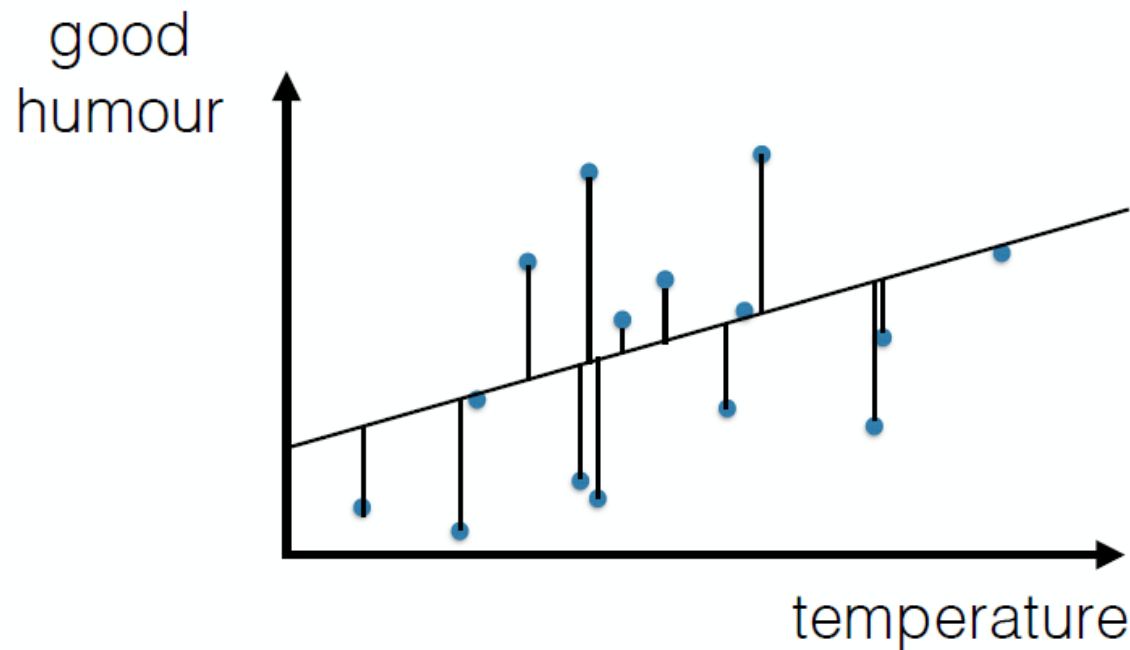


From Tanjona Ramiadantsoa

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

# Least squares

Adding covariates and  $R^2$

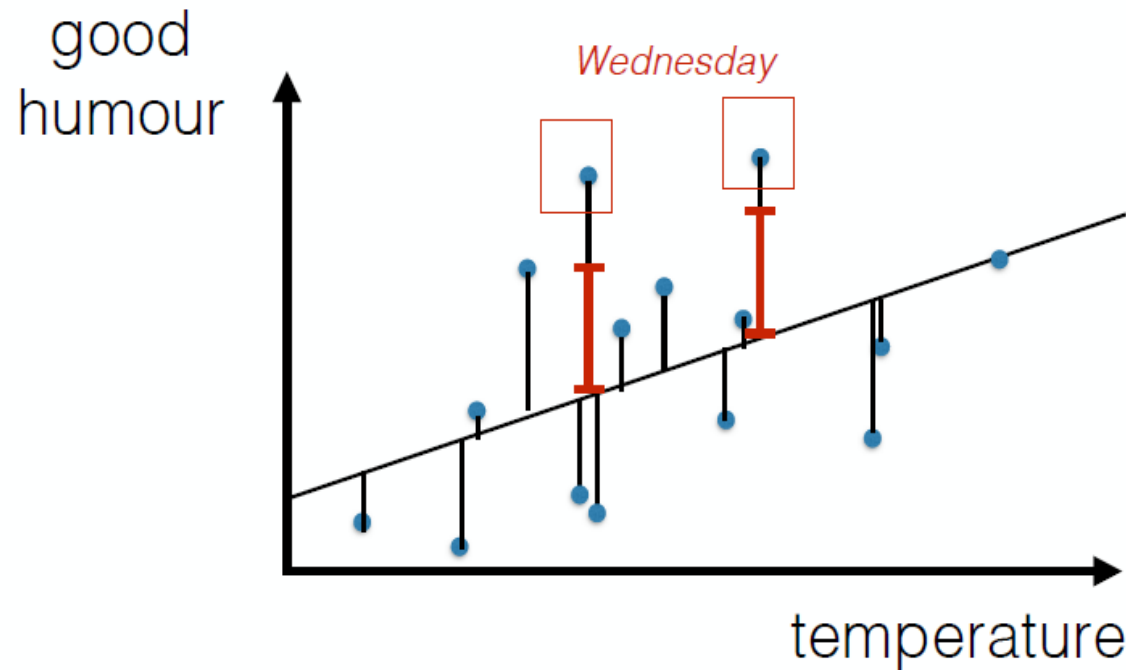


2K

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

# Least squares

Adding covariates and  $R^2$

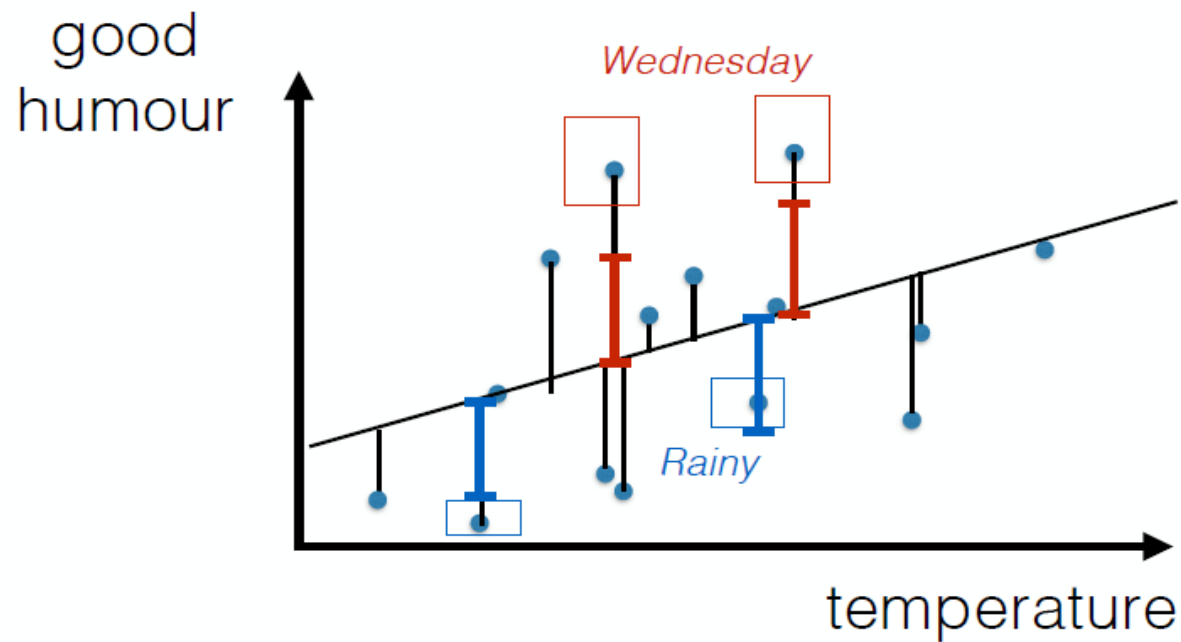


3K

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

# Least squares

Adding covariates and  $R^2$

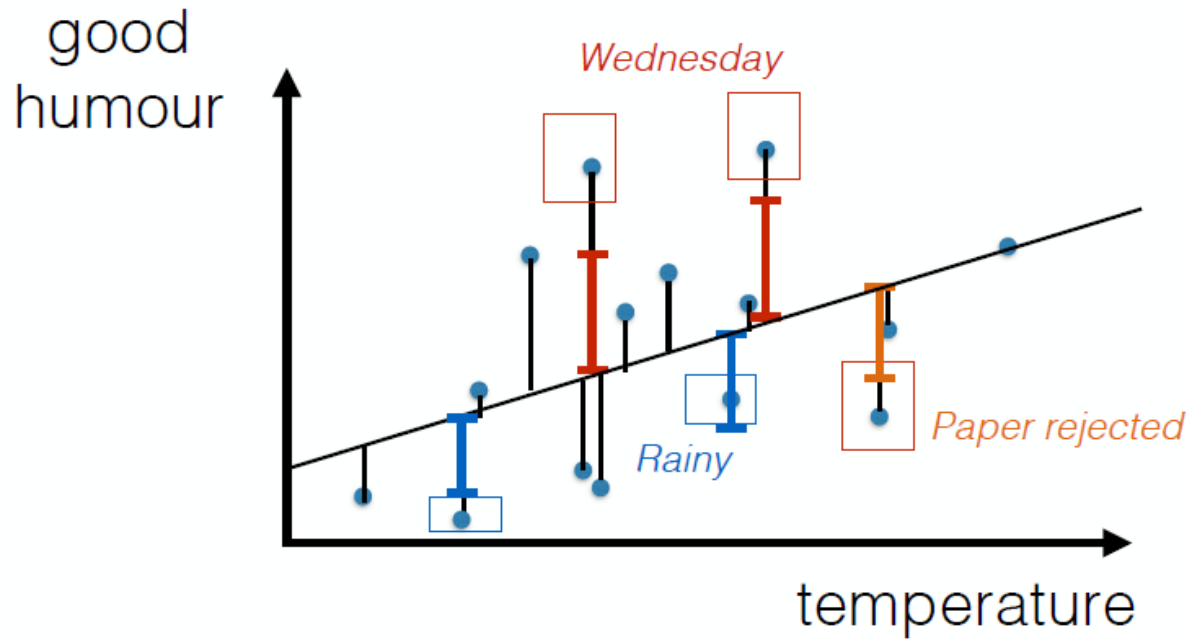


$$\text{humour} = b_0 + b_1\text{temperature} + b_2\text{Wednesday} + b_3\text{rain} + \text{Error}$$



# Least squares

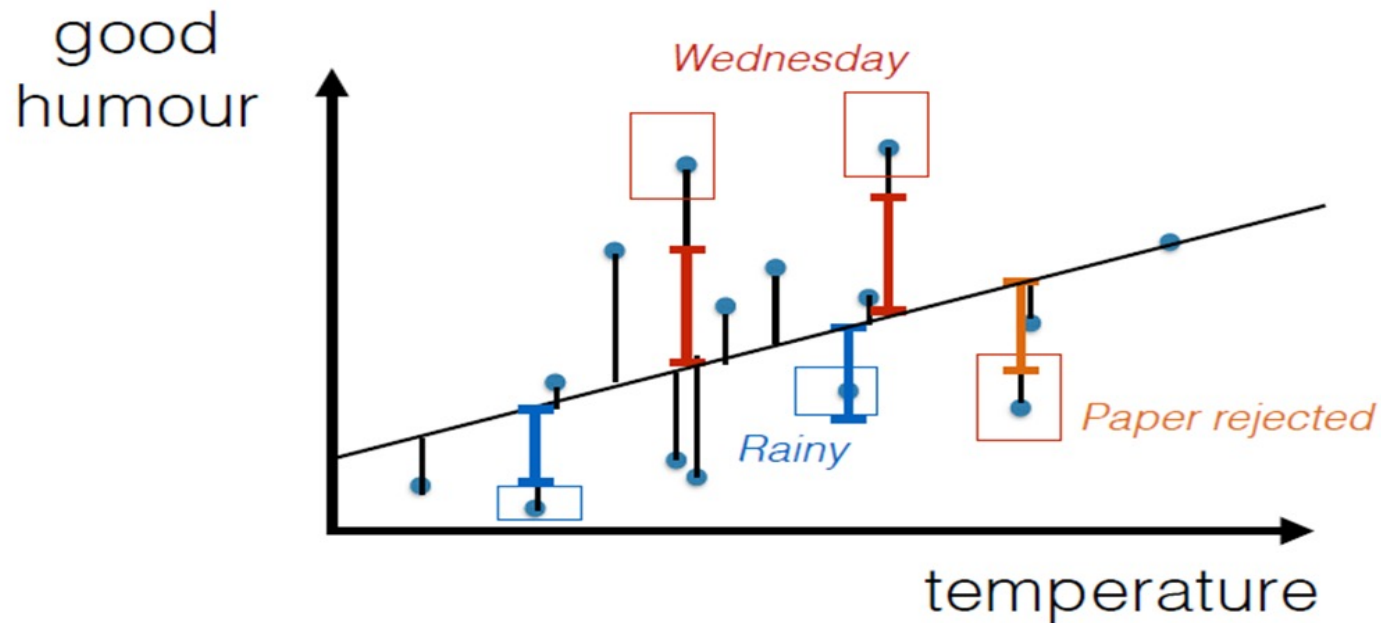
Adding covariates and  $R^2$



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

# Least squares

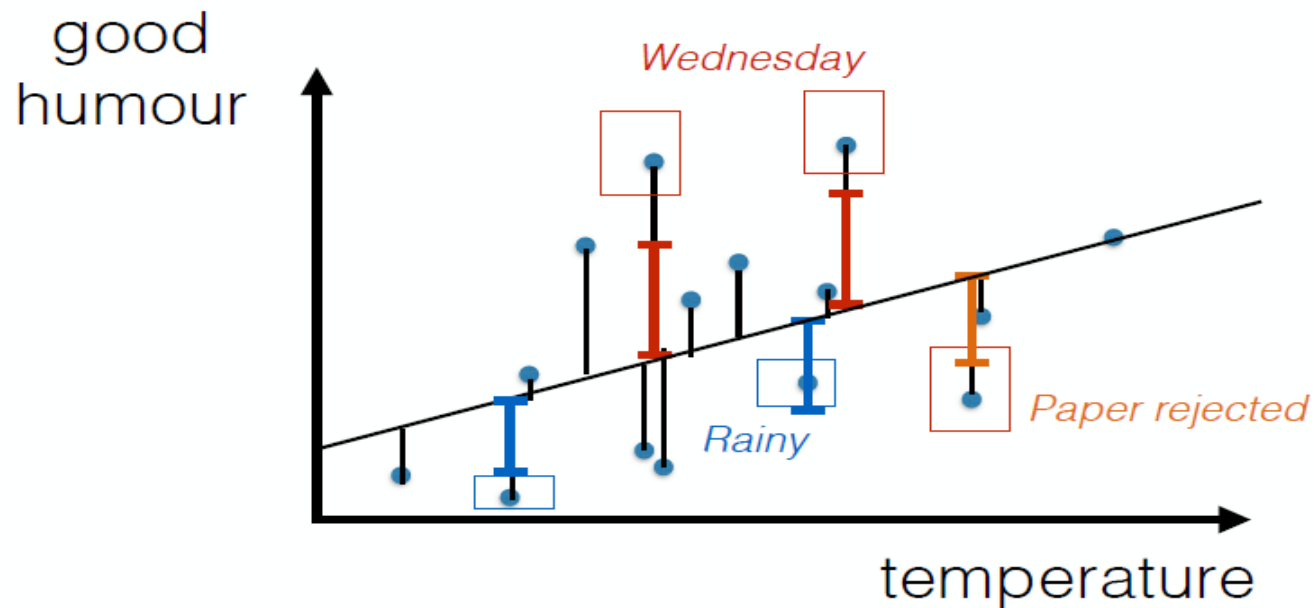
Adding covariates and  $R^2$



Adding covariates almost always increases the  $R^2$ .

# 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

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

*N*: Number of observations


*SS<sub>e</sub>*: Sum square of errors

*K*: Number of parameters

The smaller the AIC the better

# Least square AIC

More parameter is not always good

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


*N*: Number of observations

*SS<sub>e</sub>*: Sum square of errors

*K*: Number of parameters

$$(AIC = -2 \ln(L) + 2k)$$

The smaller the AIC the better

# 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](http://www.elsevier.com/locate/epidemics)



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

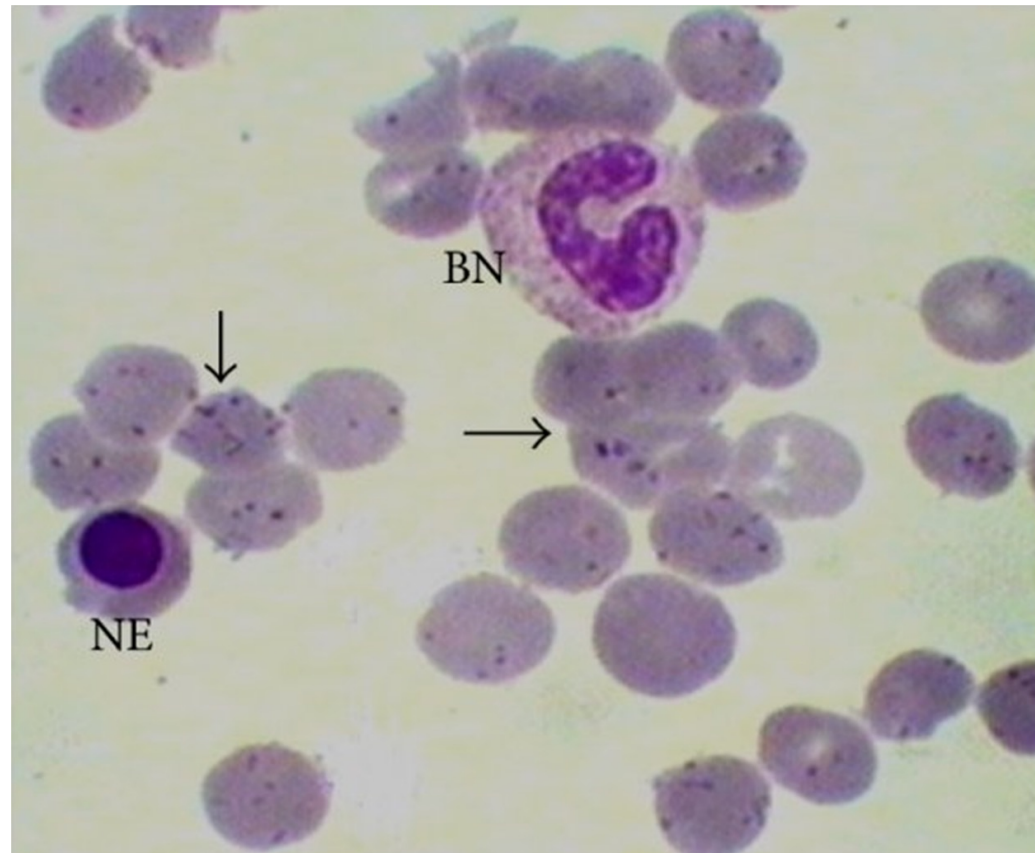


Cara E. Brook<sup>a,\*</sup>, Ying Bai<sup>b</sup>, Emily O. Yu<sup>a</sup>, Hafaliana C. Ranaivoson<sup>c,d</sup>, Haewon Shin<sup>e</sup>,  
Andrew P. Dobson<sup>a</sup>, C. Jessica E. Metcalf<sup>a,1</sup>, Michael Y. Kosoy<sup>b,1</sup>, Katharina Dittmar<sup>e,1</sup>



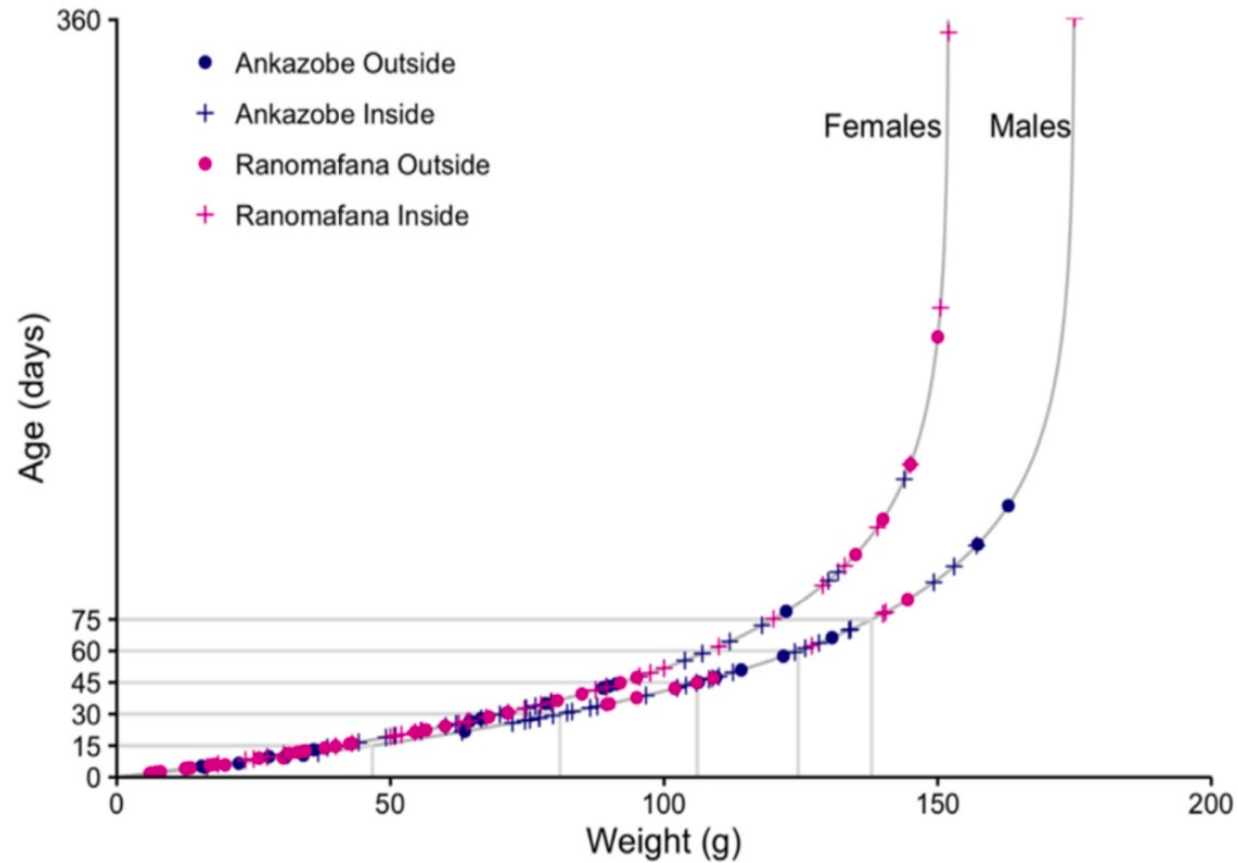
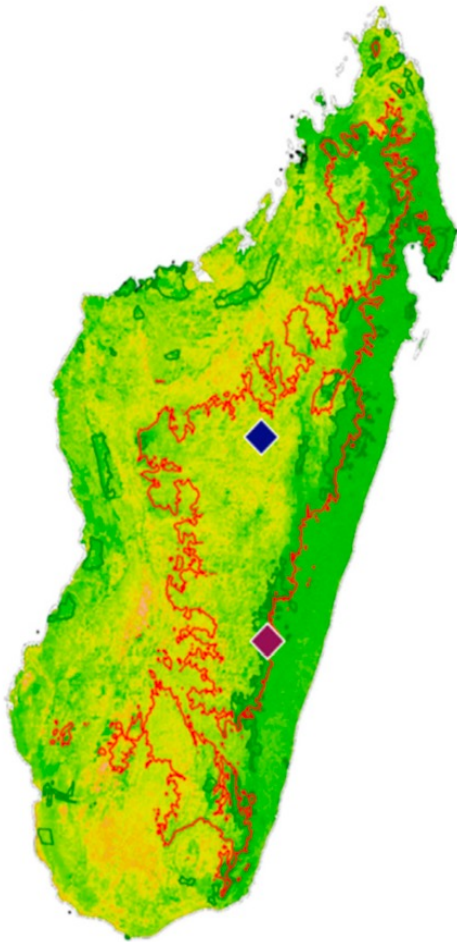
# *Bartonella* spp.

- Persistent erythrocytic bacteria that are sometimes zoonotic
- Vectored by ticks, fleas, sand flies, mosquitoes
- Some species infect humans
  - *Bartonella bacilliformis* = Carrion's disease
  - *Bartonella henselae* = cat scratch fever
  - *Bartonella quintana* = trench fever



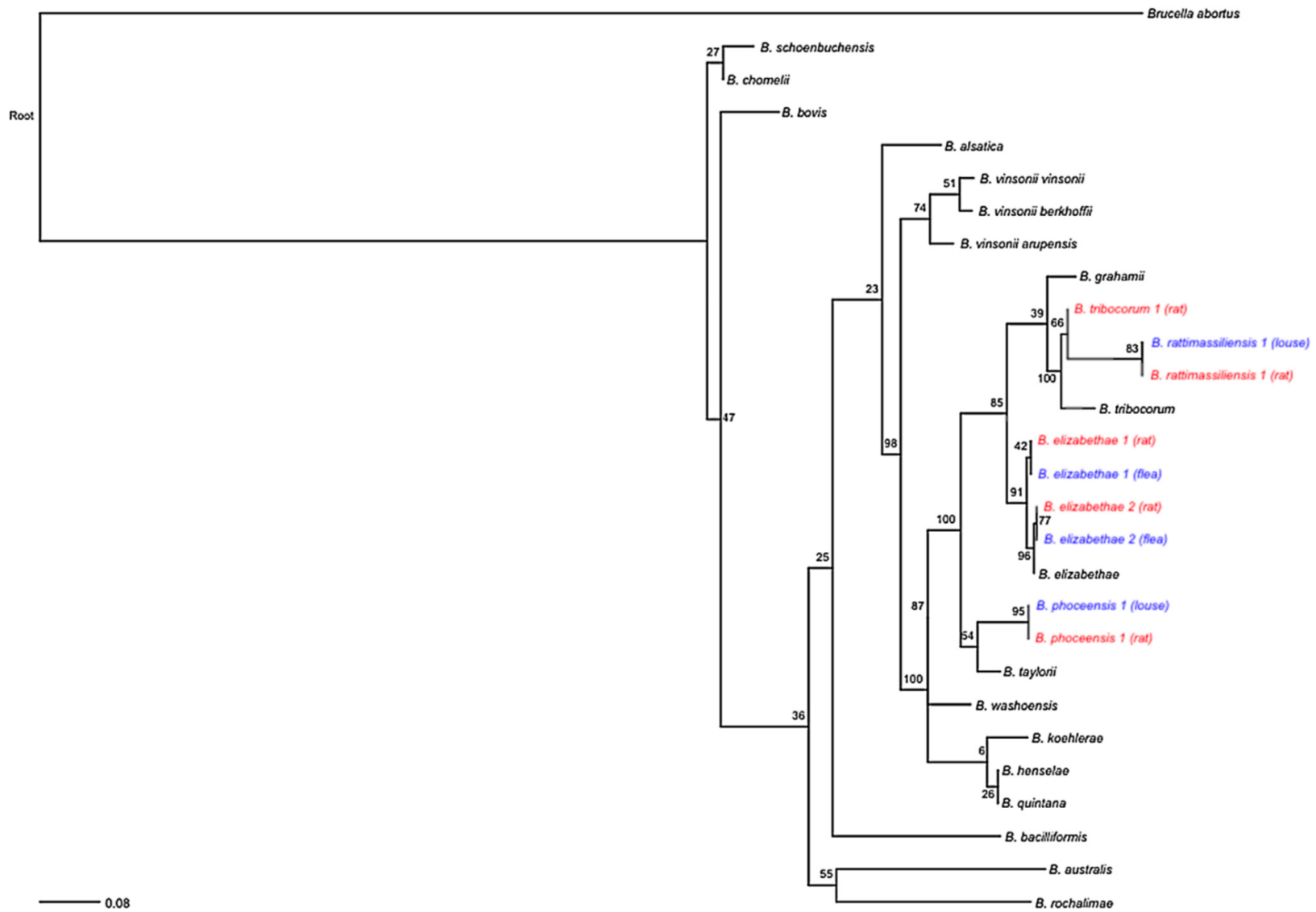


We first collected samples from rats from two sites in Madagascar.



(Ricker 1979)

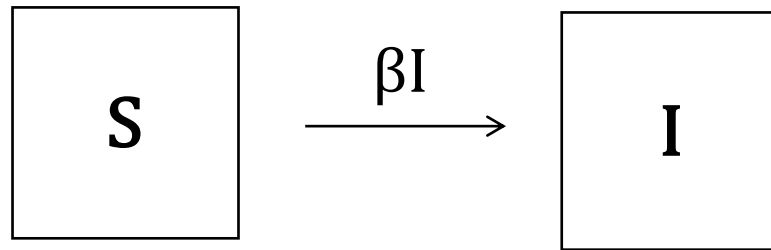
Statistically, we demonstrated an association between genotypes of *Bartonella* spp. found in rats and their ectoparasites.



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

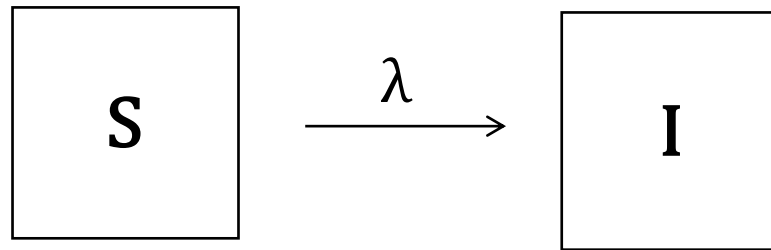
age	sick	age.class	sex
1.685375	0	1	M
2.011821	0	1	F
2.344301	0	1	F
2.611943	0	1	F
4.145697	0	1	F
4.319159	0	1	F
4.319159	0	1	F
4.493246	0	1	F
4.662362	0	1	M
4.784125	0	1	M

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



for a persistent, non-immunizing infection

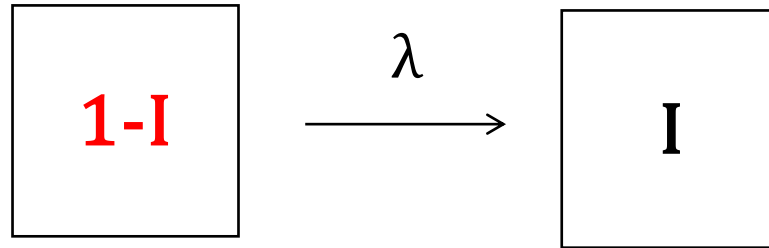
Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.



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

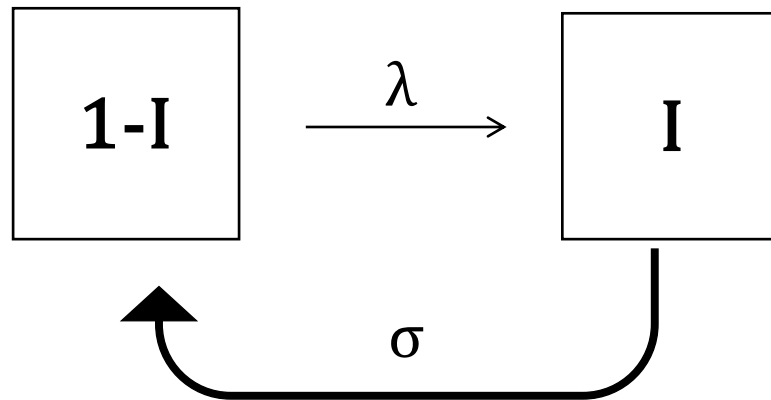
Age-prevalence data allows for powerful inference into the dynamics of pathogen transmission.

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



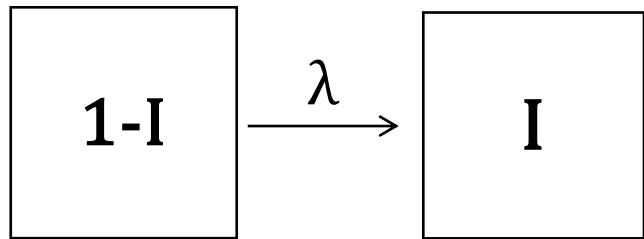
where  $\lambda$ , the force of infection, is the per capita rate at which susceptible hosts become infected

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



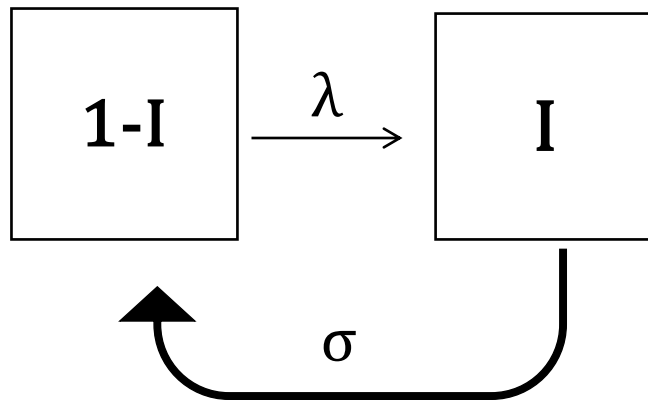
and  $\sigma$  is the rate of recovery from infection

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



A diagram showing a transition from a box labeled  $1-I$  to a box labeled  $I$ . An arrow points from the first box to the second, with the Greek letter  $\lambda$  written above it.

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

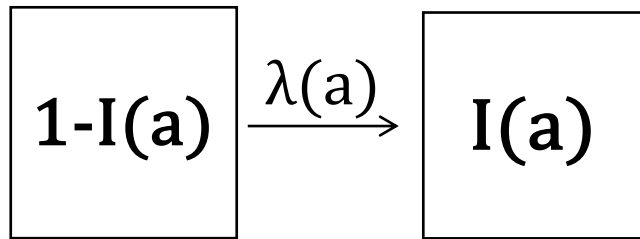


A diagram showing two boxes, one labeled  $1-I$  on the left and one labeled  $I$  on the right. An arrow points from the left box to the right box, with the Greek letter  $\lambda$  written above it. A curved arrow points from the bottom of the right box back to the bottom of the left box, with the Greek letter  $\sigma$  written below it.

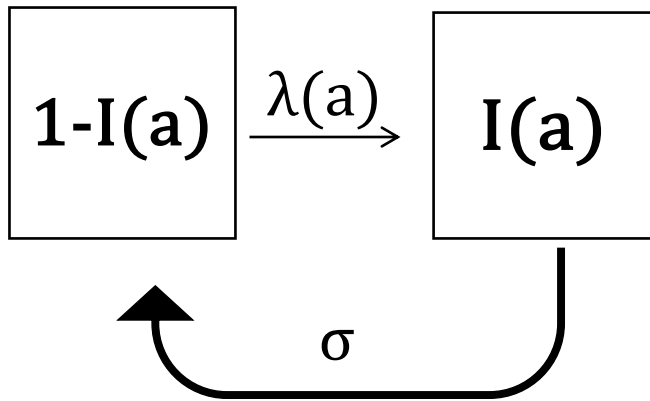
$$\frac{dI}{dt} = \lambda(1-I) - \sigma I$$



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

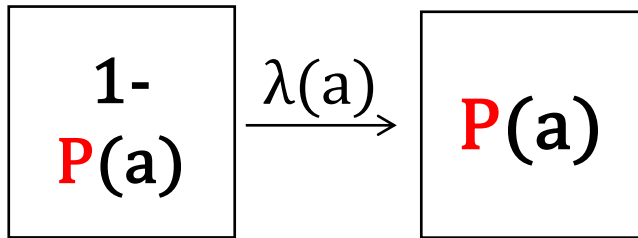


$$\frac{dI(a)}{da} = \lambda(a)(1 - I(a))$$

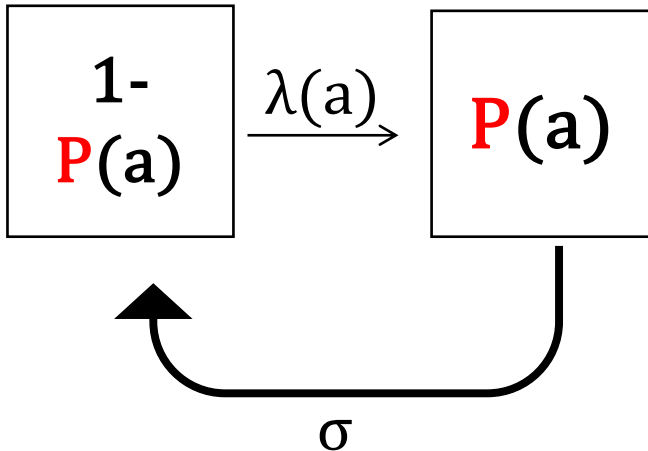


$$\frac{dI(a)}{da} = \lambda(a)(1 - I(a)) - \sigma I(a)$$

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



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



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

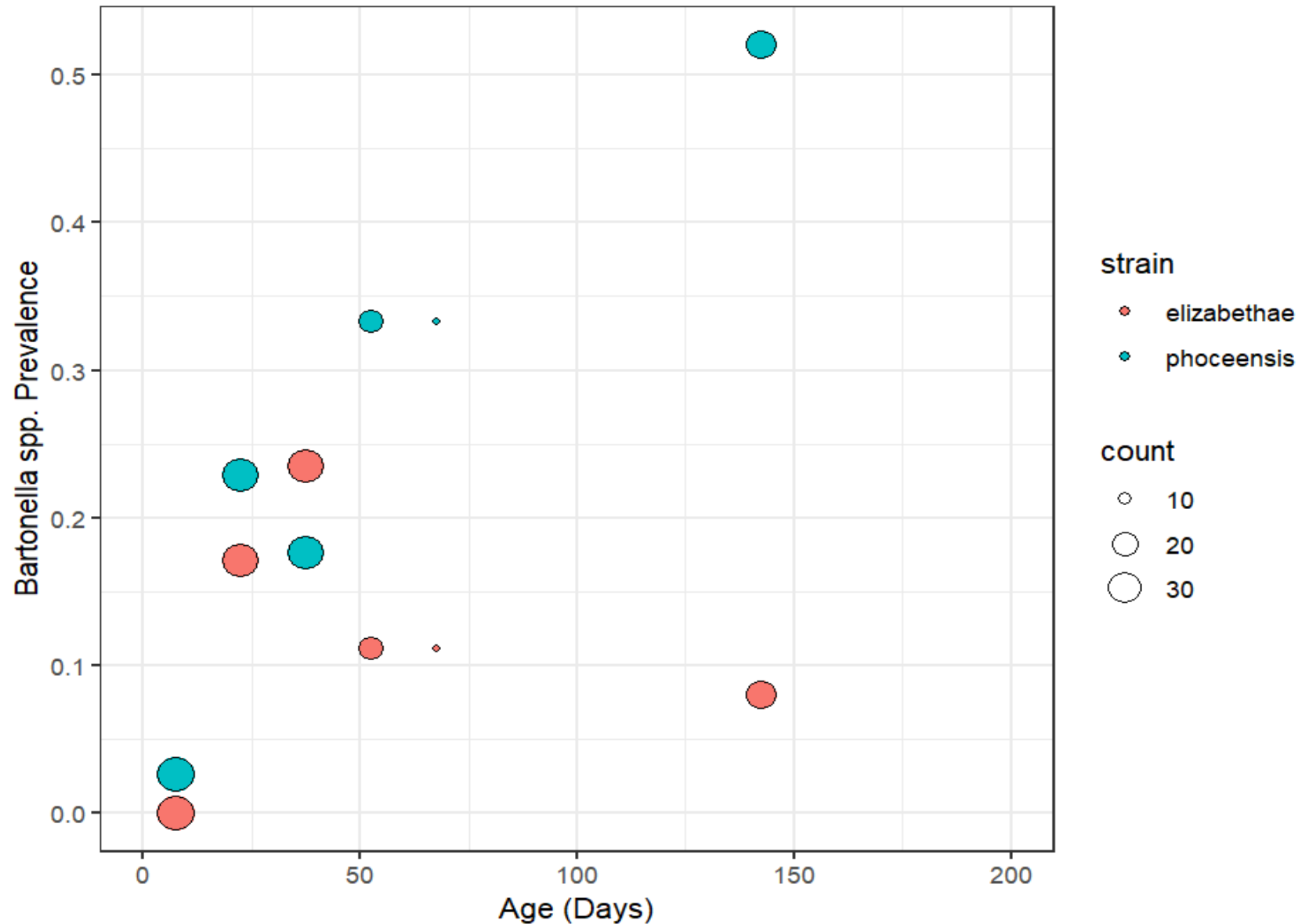
Compare using  **$AIC = 2K - 2\ln(L)$**

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

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

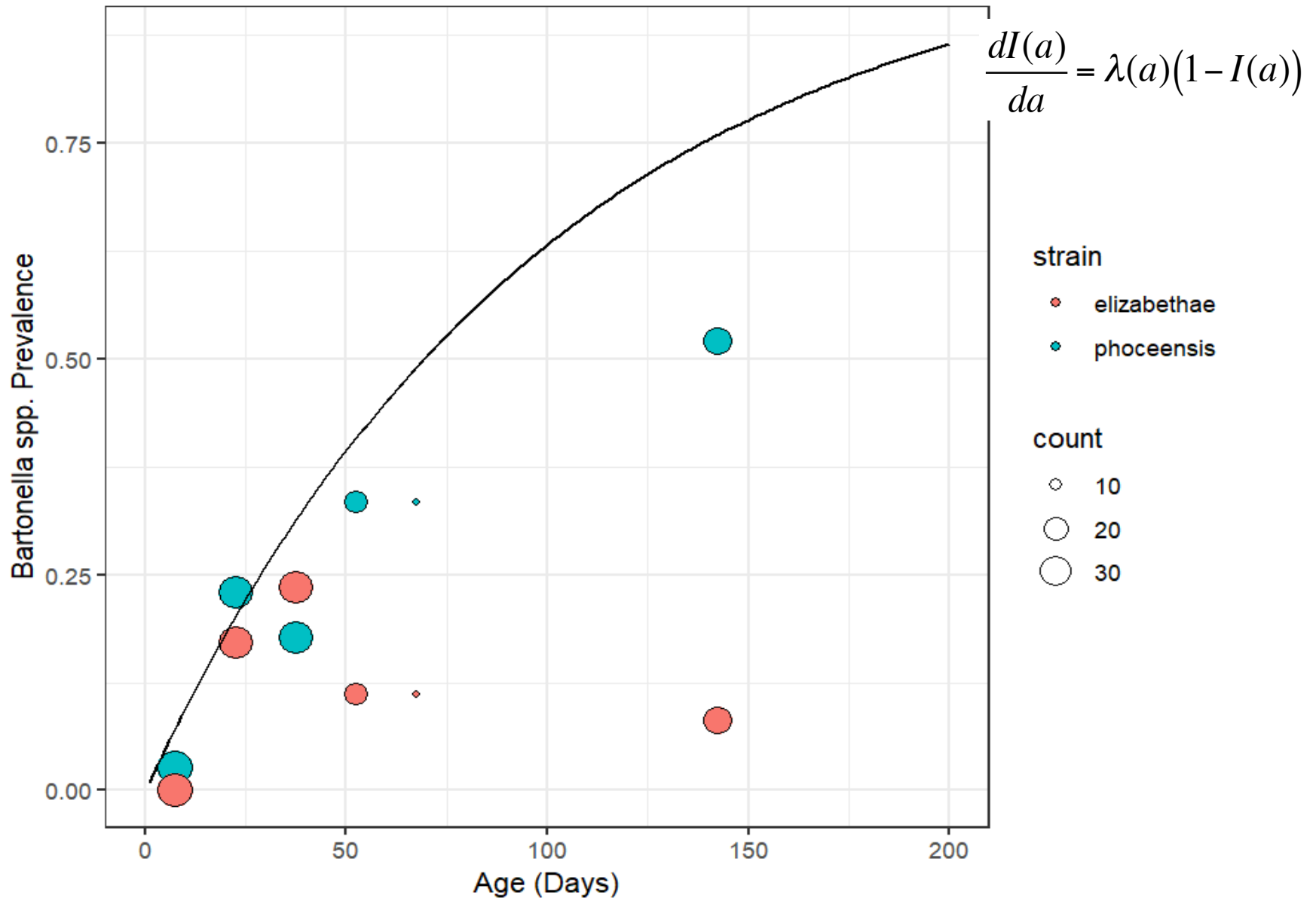
# Look at the data !

Jereo aloha hoe manao ahoana



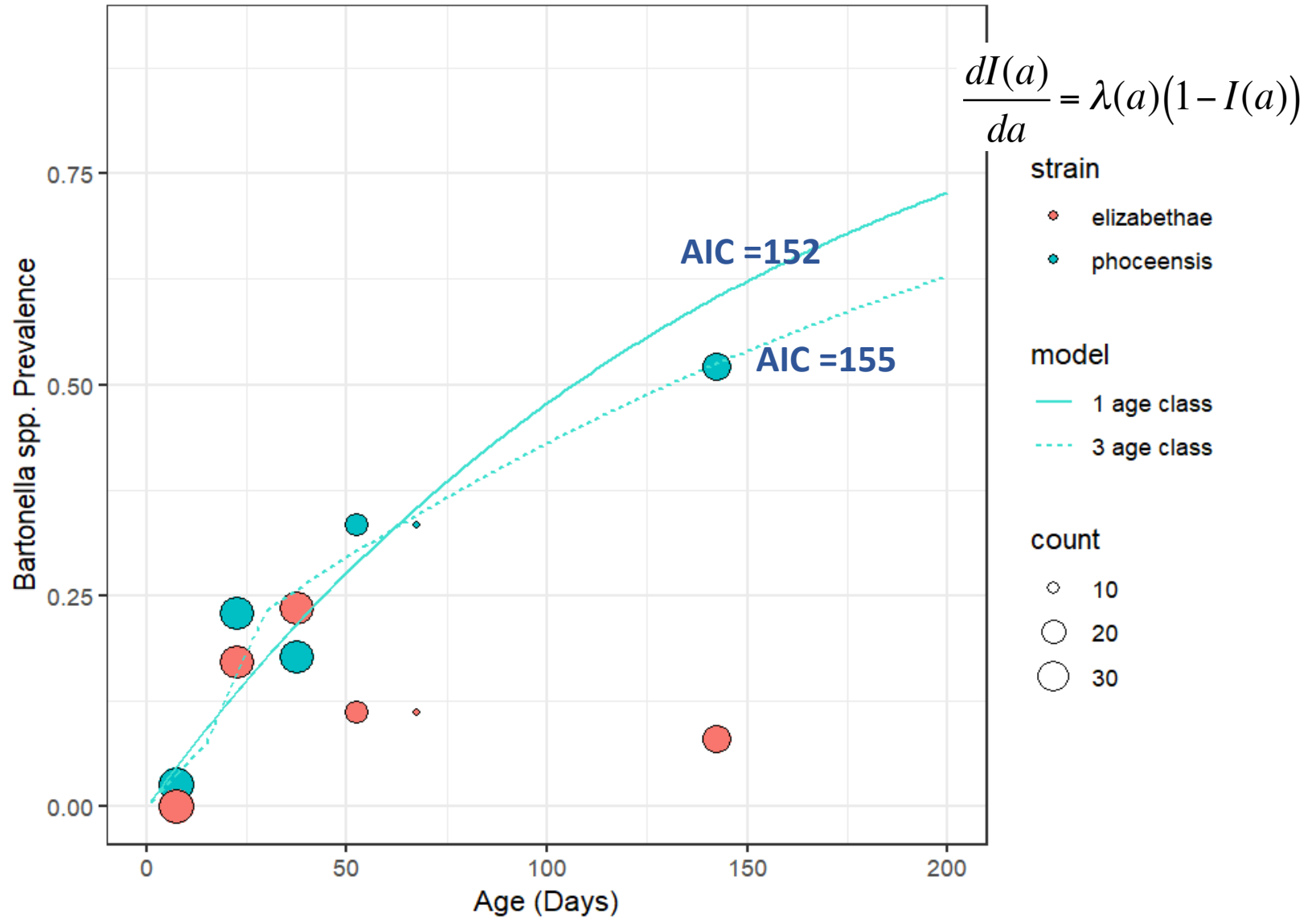
# Try the model!

Andramo kely

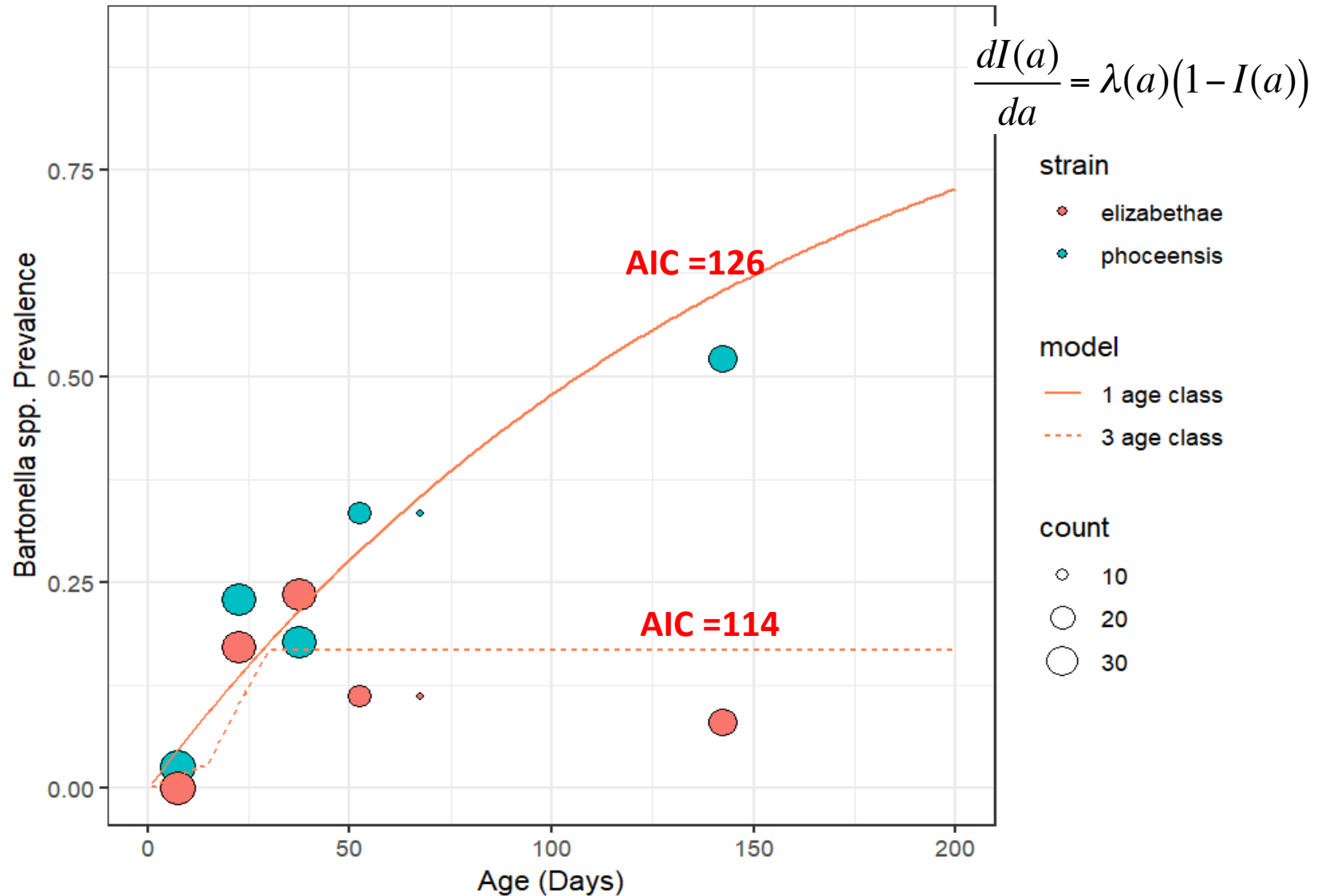


# Keep trying!

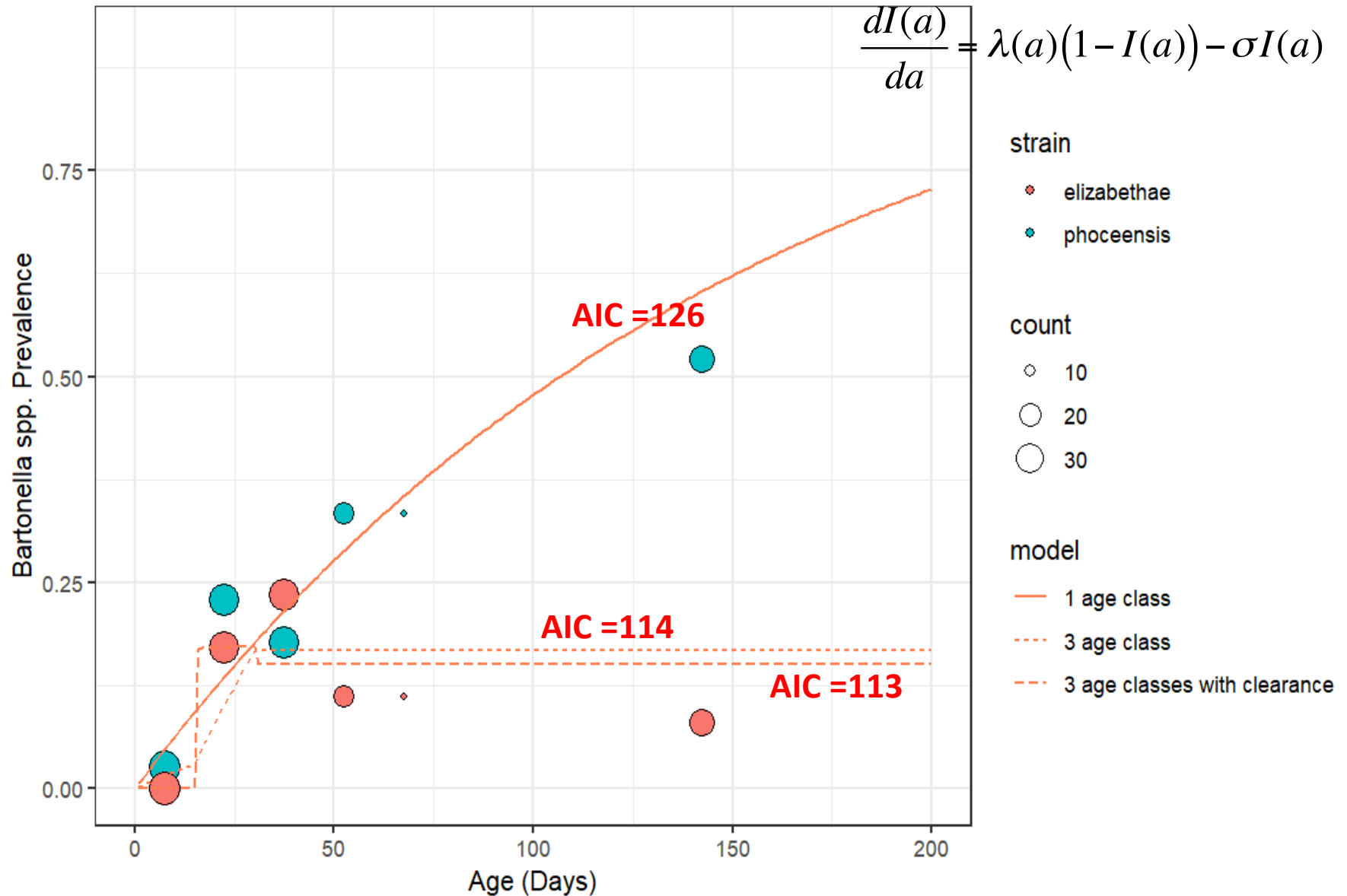
# Alô fô !



# Keep trying!

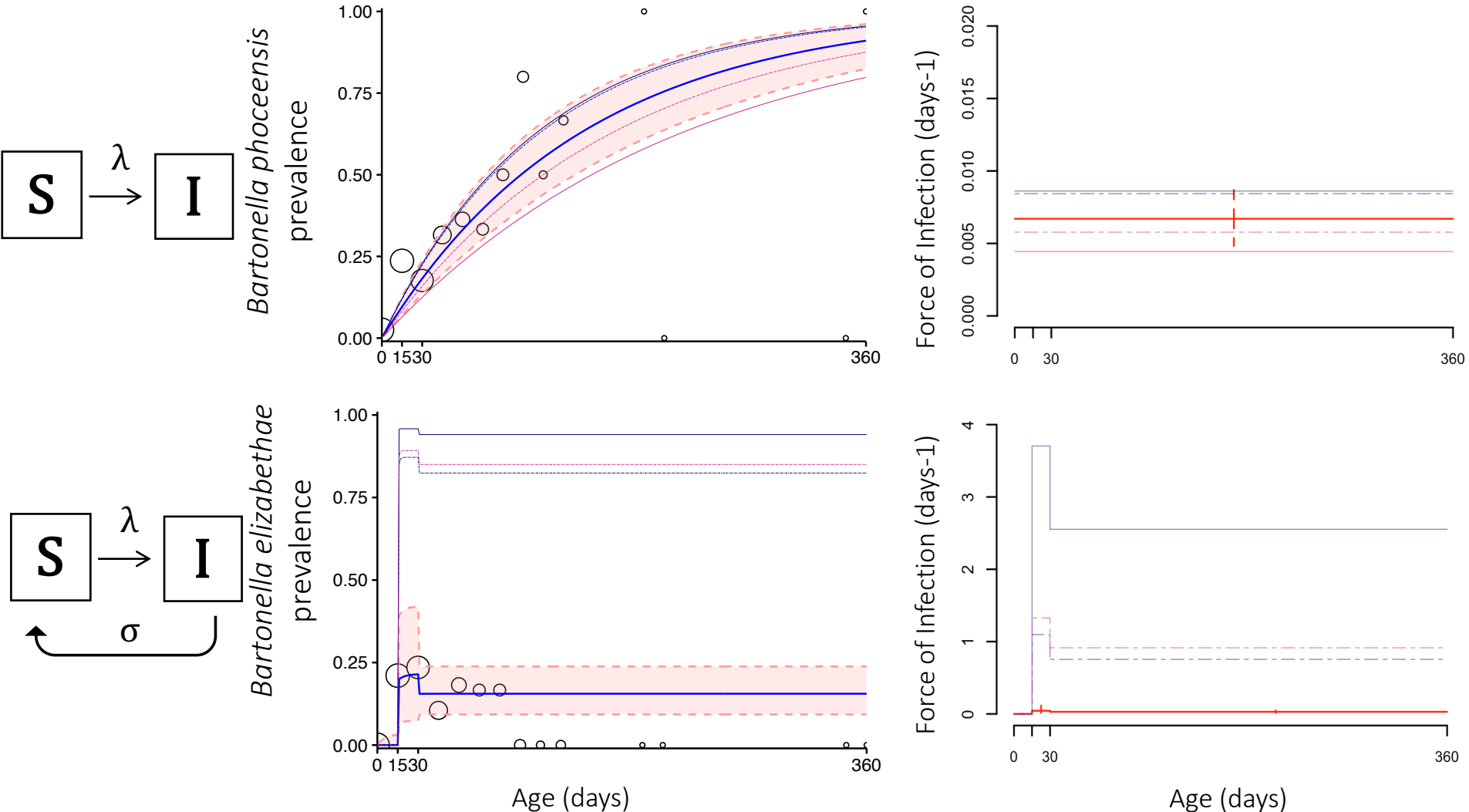


# Keep trying!





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.

