Model Selection and Comparison



Christian Ranaivoson E2M2 2025 Andasibe, Madagascar

University of Chicago, USA ; Association Ekipa Fanihy; University of Antananarivo, Madagascar Adapted from Cara Brook, Jessica Metcalf, Christian Ranaivoson E2M2 2022-2024.

- Introduction to model selection

- Be familiar with common measures of model fit
- Tutorial on how to do it in practice (Real example)

Which model is best?



Date of symptoms onset

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 data and your question.

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 data and your question.

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

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

Maximum likelihood (Maximum de vraisemblance)

(izay manakaiky indrindra ny tena izy ry reto an)

R-squared (R-carré)

Least squares (Moindres carrés)

AIC

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



Maximum likelihood



e.g., height

Maximum likelihood







Likelihood



~ the 'probability of seeing the data' given the chosen parameters ...each individual in the data has a corresponding likelihood... multiply them...

Very low likelihood.

Maximum likelihood











~ the 'probability of seeing the data' given the chosen parameters ...each individual in the data has a corresponding likelihood... multiply them...

Very low likelihood.



Likelihood Vraisemblance

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

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

Examples:

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

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





Likelihood Vraisemblance

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

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

Examples:

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

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





Optimization/maximization

A function and its derivative

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



From Tanjona Ramiadantsoa

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 find minimum/maximum.

Michael Jordan's FT success in 100 match career







Definition r²



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



Adding covariates and R²



humour = $b_0 + b_1$ temperature + b_2 Wednesday +Error

Adding covariates and R^2



humour = $b_0 + b_1$ temperature + b_2 Wednesday+

b₃rain +Error

Adding covariates and R²



humour = $b_0 + b_1$ temperature + b_2 Wednesday+ b_3 rain + b_4 rejection +Error

Adding covariates and R²



Adding covariates almost always increases the R²

Adding covariates and R²



Adding covariates almost always increases the R² - so a key question is when to stop.

What to choose?



Least square AIC

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

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

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

N: Number of observations SS_e: 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



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
- 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?

sick [‡]	age.class 🎈	sex 🍦
0	1	М
0	1	F
0	1	F
0	1	F
0	1	F
0	1	F
0	1	F
0	1	F
0	1	М
0	1	М
	sick	sick











Simple SI model

Ν







for a persistent, non-immunizing infection



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

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



and $\boldsymbol{\sigma}$ is the rate of recovery from infection



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



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



similar techniques can also be applied to ageseroprevalence 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.



