
VEROSSIMILHANÇA MÁXIMA E SELEÇÃO DE MODELOS



Hipótese (θ_1)

Dado honesto

Qual a probabilidade de arremessarmos um dado **5 vezes** e obtermos **6 em todos os arremessos (x)**?



Hipótese (θ_1)

Dado honesto

$$P(\mathbf{x}|\theta_1) = P(P_1 * P_2 * \dots P_5)$$
$$(1/6)^5 = 1/7776$$



Hipótese (θ_1)

Dado honesto

$$P(x|\theta_1) = 1/7776 = 0.00012$$





Hipótese (θ_1)

Dado honesto

$$P(\mathbf{x}|\theta_1) = 1/7776 = 0.00012$$

Hipótese (θ_2)

Dado viciado!!

+ VEROSSÍMIL

$$P(\mathbf{x}|\theta_2) = (0.8)^5 = 0.33$$

$$L(\theta|\mathbf{x}) = P(\mathbf{x}|\theta)$$

A **verossimilhança** (**L**) de um conjunto de **parâmetros** (**θ**), dada uma **observação** (**\mathbf{x}**) é igual a probabilidade desta mesma observação ter ocorrido dado os valores daqueles parâmetros.

$$L(\theta|x) = P(x|\theta)$$

$$\neq P(\theta|x)!!$$

A **verossimilhança** (L) de um conjunto de **parâmetros** (θ), dada uma **observação** (x) é igual a probabilidade desta mesma observação ter ocorrido dado os valores daqueles parâmetros.

Observação (x): Barulhos vindos do forro da sua casa.

Hipótese (θ_1): *Gremlins* estão dando uma festa!



Observação (x): Barulhos vindos do forro da sua casa.

Hipótese (θ_1): *Gremlins* estão dando uma festa!



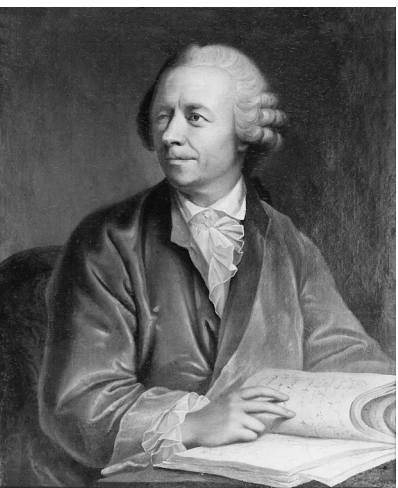
$P(x|\theta_1)$ é altamente **provável e verossímil** $L(\theta_1|x) \dots$

\dots mas $P(\theta_1|x)$ é altamente **improvável**



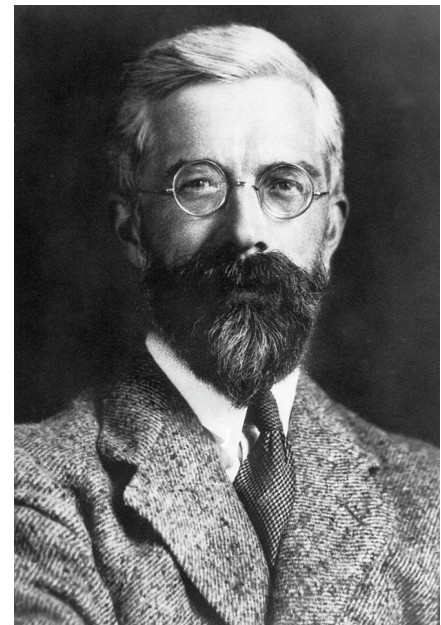
Daniel Bernoulli (1700-1782)

$$L(\theta|\mathbf{x}) = P(\mathbf{x}|\theta)$$

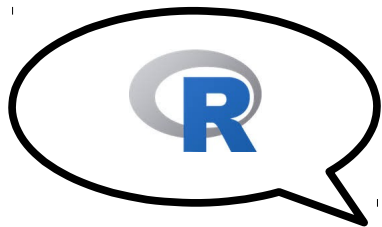


Leonhard Euler (1707-1783)

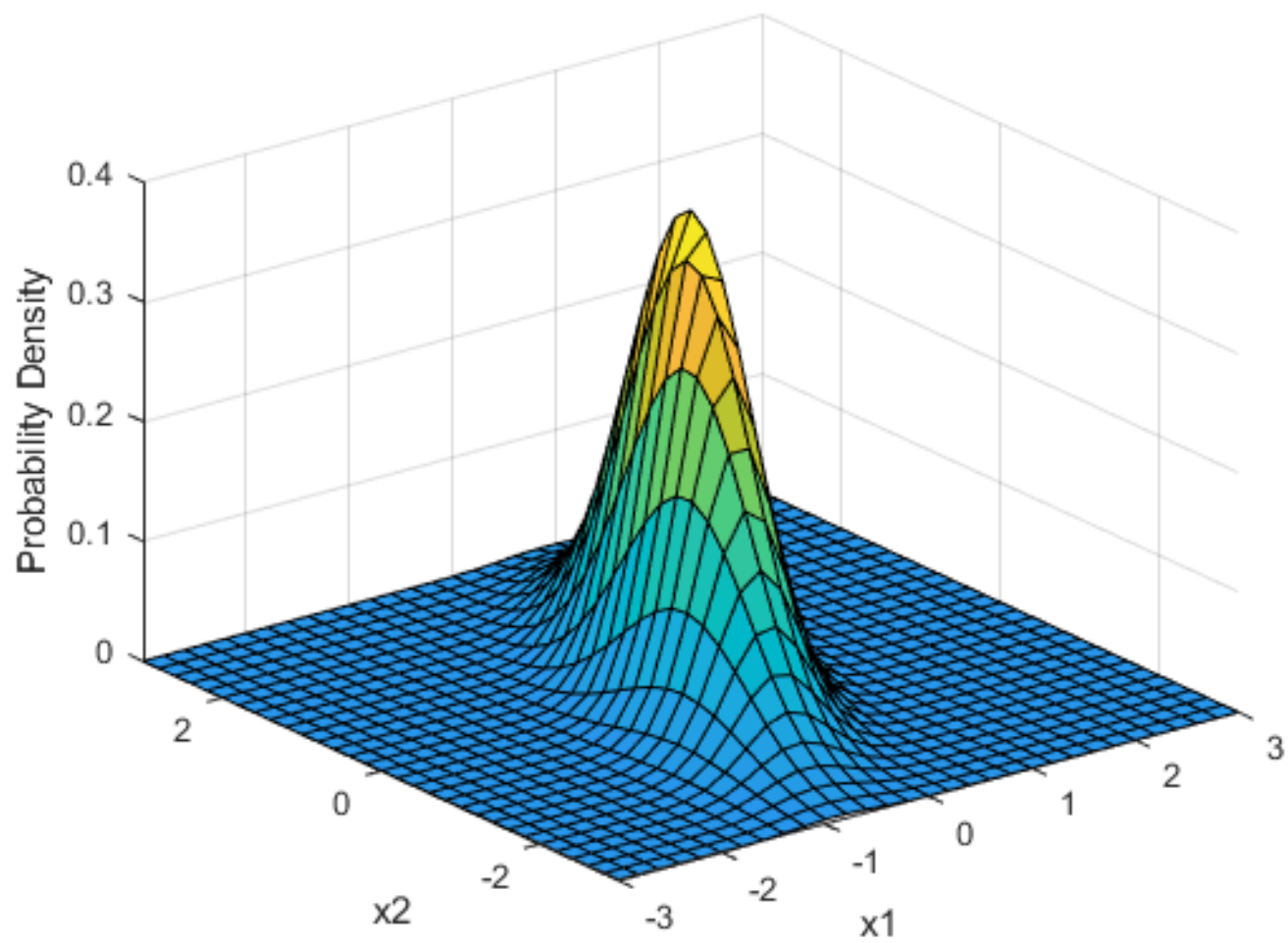
A verossimilhança pode ser usada para **estimar** um ou mais parâmetros (θ)



Ronald Fisher (1890-1962)



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

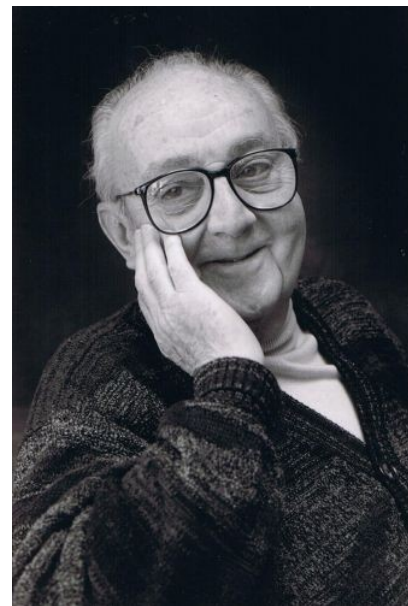


Então... o modelo com **mais parâmetros** (+ ajustados) é sempre o **mais verossímil**?



Mas afinal... o que são **modelos**?

*“Essentially,
all models are wrong,
but some are useful.”*



George Box (1919-2013)

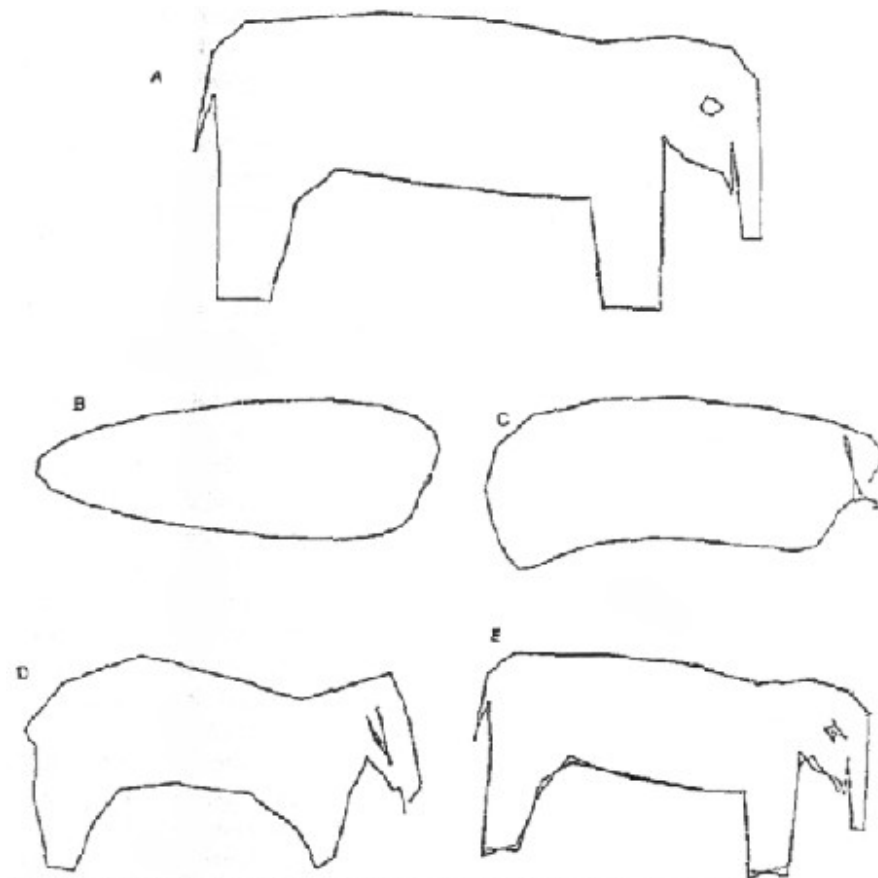


FIGURE 1.2. "How many parameters does it take to fit an elephant?" was answered by Wei (1975). He started with an idealized drawing (A) defined by 36 points and used least squares Fourier sine series fits of the form $x(t) = \alpha_0 + \sum \alpha_i \sin(i t \pi / 36)$ and $y(t) = \beta_0 + \sum \beta_i \sin(i t \pi / 36)$ for $i = 1, \dots, N$. He examined fits for $K = 5, 10, 20$, and 30 (shown in B-E) and stopped with the fit of a 30 term model. He concluded that the 30-term model "may not satisfy the third-grade art teacher, but would carry most chemical engineers into preliminary design."

A navalha de Occam

e o Princípio da parsimônia

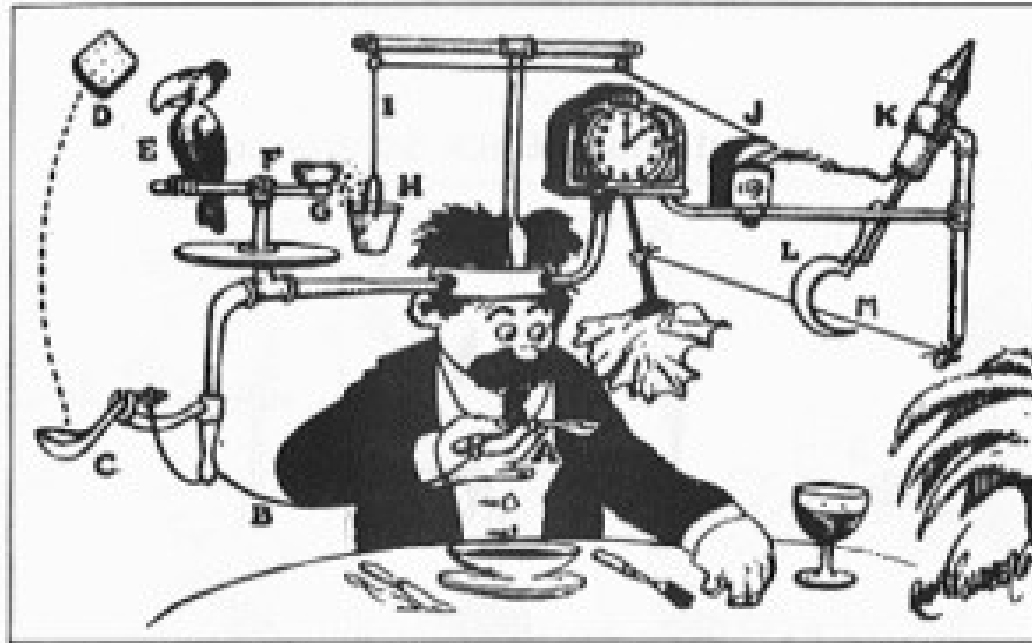
*A pluralidade não
deve ser posta
sem necessidade*

*When you have two competing
theories that make exactly the
same predictions, the simpler one
is the better.*



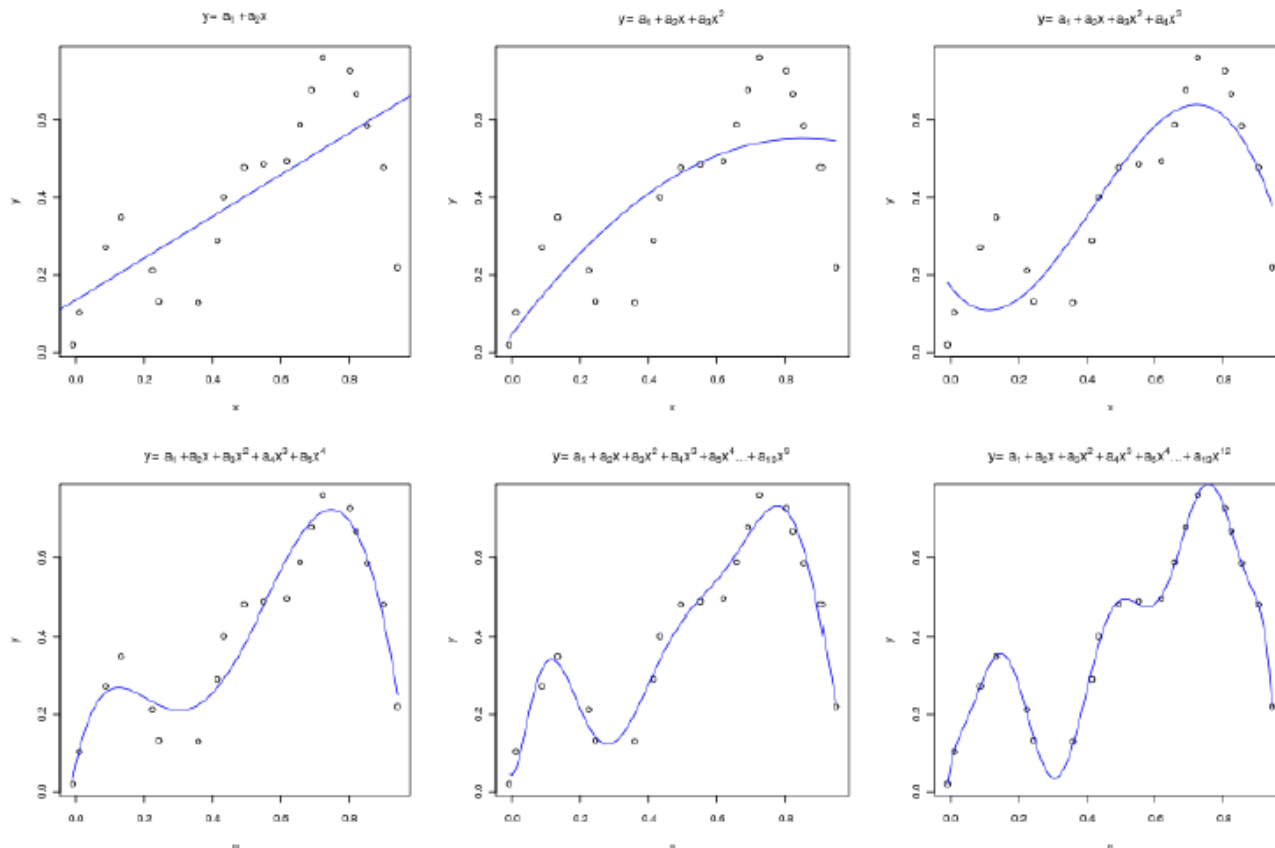
William de Occam (1285-1347)

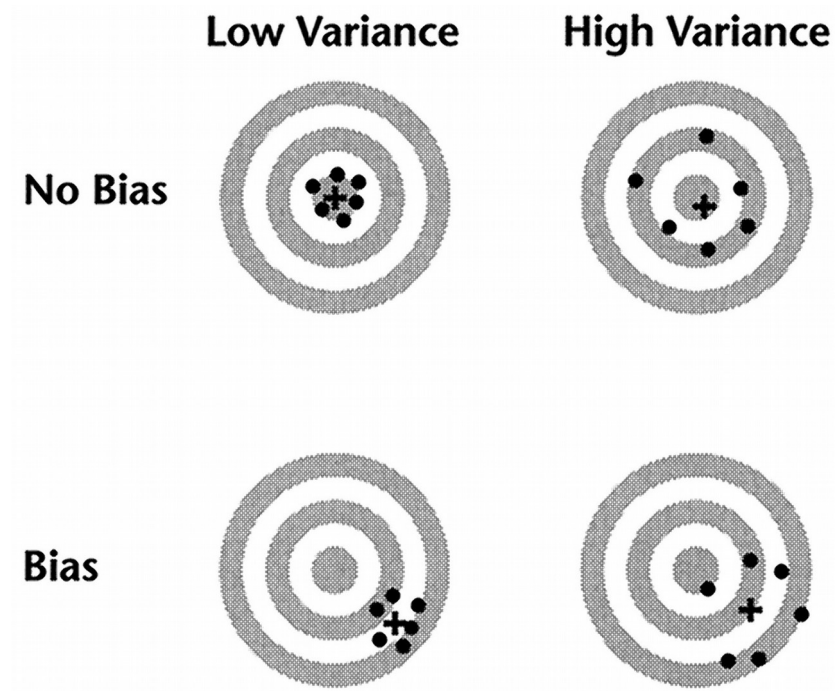
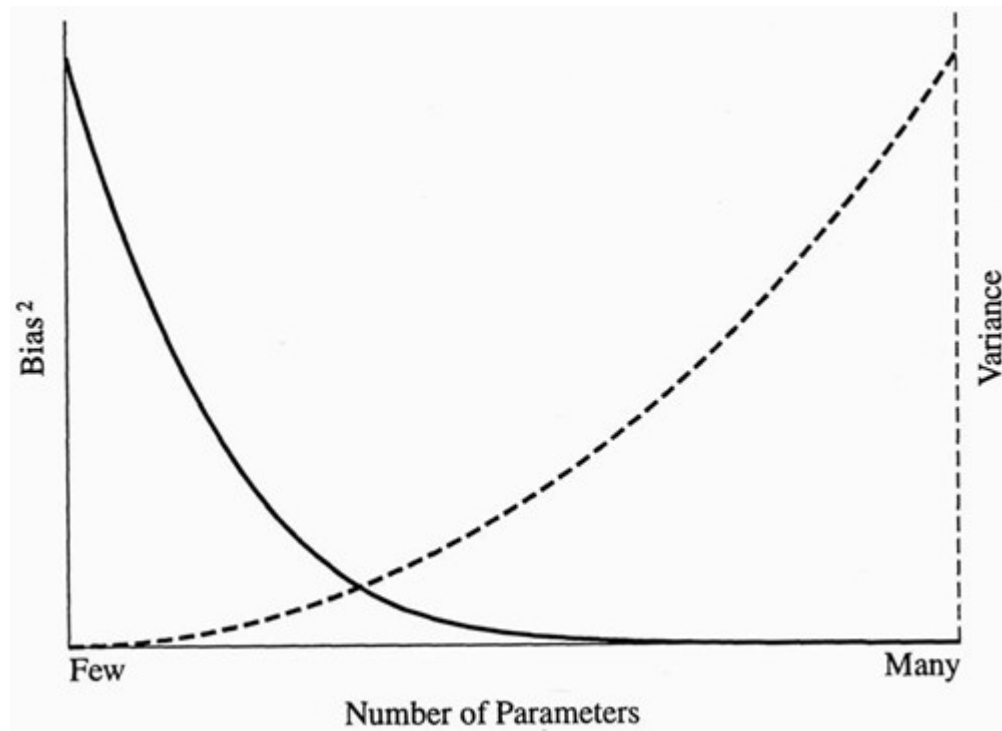
Self-Operating Napkin



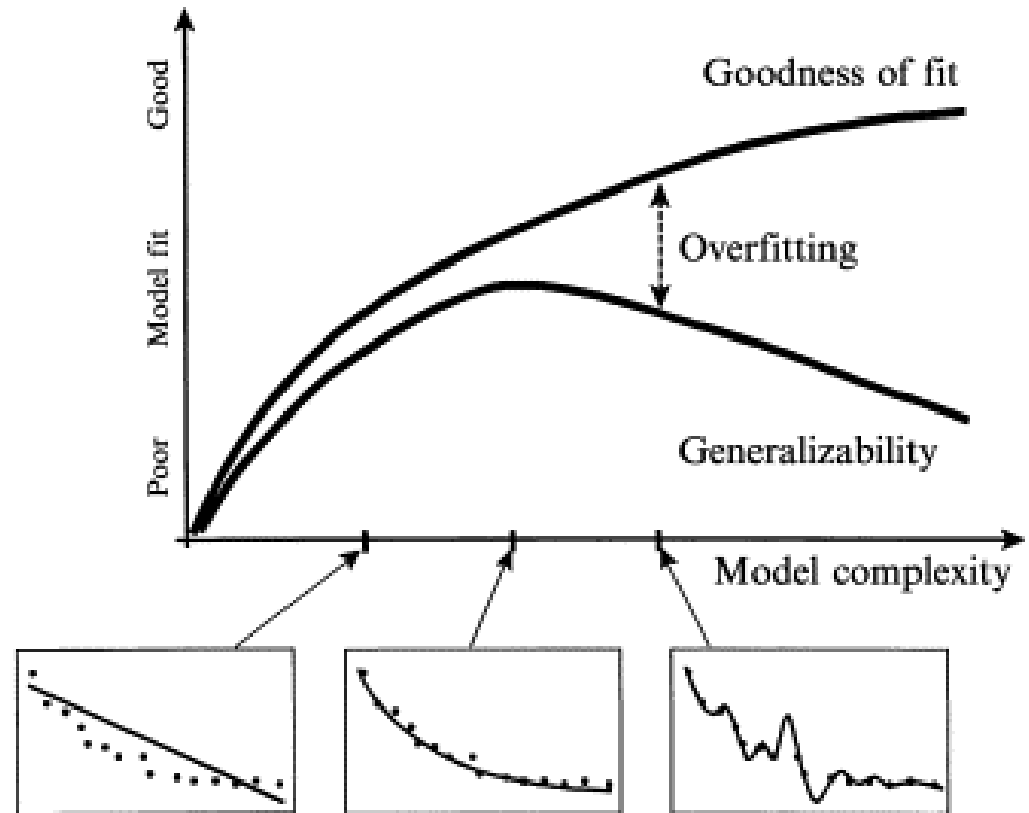
Rube Goldberg

Sobreajuste (Overfitting)





Na ecologia,
não queremos **apenas**
explicar mas também
prever.



Dragagem de dados *(data dredging)*

“Data dredging is the use of **data mining** to uncover patterns in data that can be presented as statistically significant without first devising a specific hypothesis as to the underlying causality”

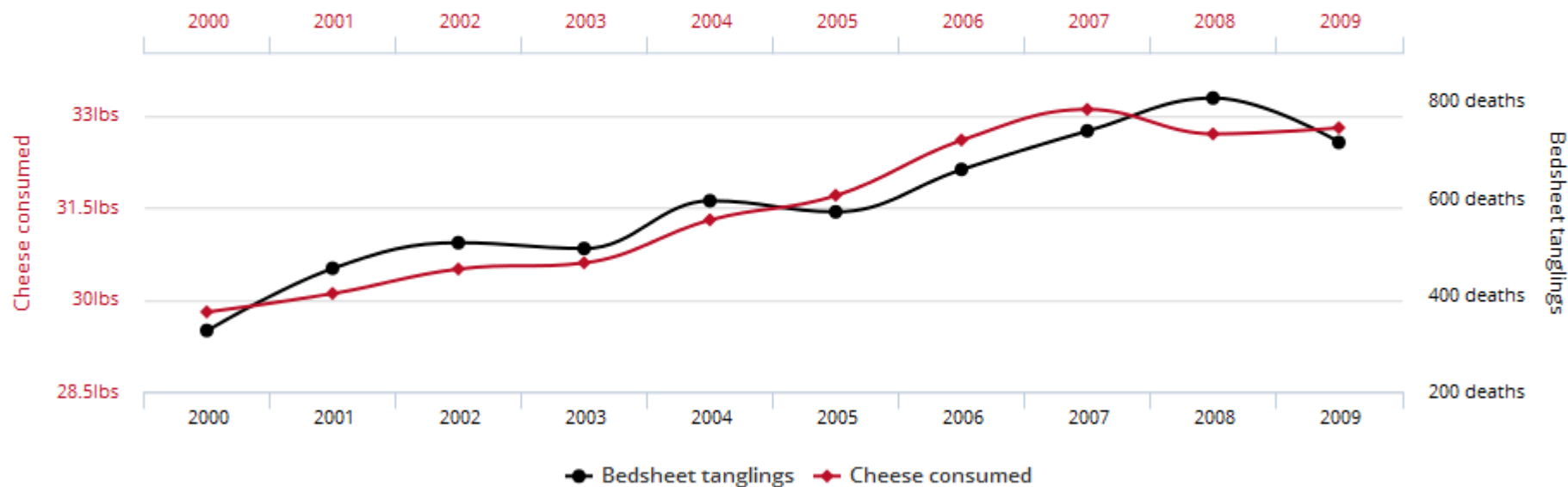


Per capita cheese consumption

correlates with

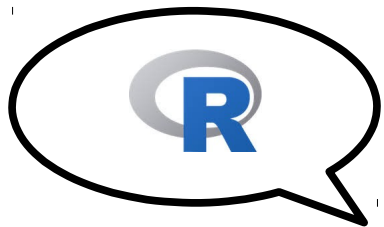
Number of people who died by becoming tangled in their bedsheets

Correlation: 94.71% ($r=0.947091$)



Data sources: U.S. Department of Agriculture and Centers for Disease Control & Prevention

tylervigen.com



Teoria da Informação + Verossimilhança

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



Hirotugu Akaike (1927-2009)

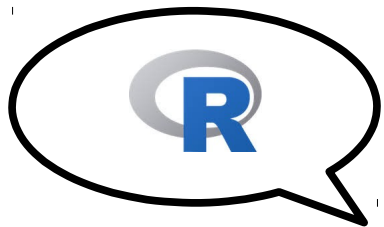
Modelo	k	ln(L)	AIC	Δ_i	ω_i
x_1	β_0, β_1	-27	58	0	0.6224
$x_1 + x_2$	$\beta_0, \beta_1, \beta_3$	-26,5	59	1	0.3775
x_2	β_0, β_1	-78	92	34	0.0001

$$\Delta_i = AIC_i - AIC_{min}$$

$$w_i = \frac{\exp\left(-\frac{1}{2} \Delta_i\right)}{\sum_{r=1}^R \exp\left(-\frac{1}{2} \Delta_r\right)}$$



Hirotugu Akaike (1927-2009)



Checklist

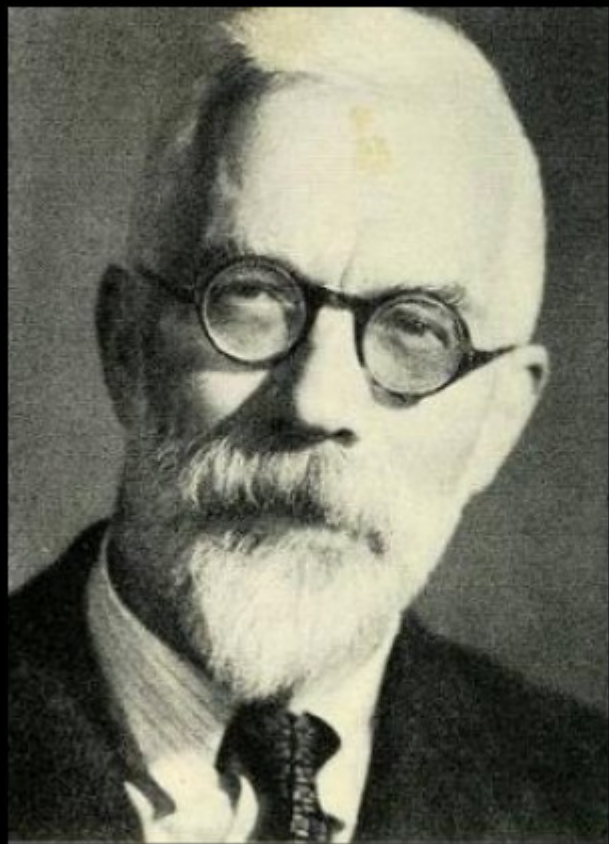


1. Defina claramente a sua **pergunta de estudo**;
2. Selecione seu conjunto de modelos/covariáveis) com base no seu **conhecimento biológico/ecológico**;
3. Seja **parcimonioso(a)**;

Checklist



1. NÃO compare amostras/variáveis dependentes diferentes;
2. NÃO compare variáveis dependentes em escalas diferentes!
3. Lembre-se de que não está testando significância ou rejeitando hipóteses.



To call in the statistician after the experiment is done may be no more than asking him to perform a post-mortem examination: he may be able to say what the experiment died of.

- Ronald Fisher -