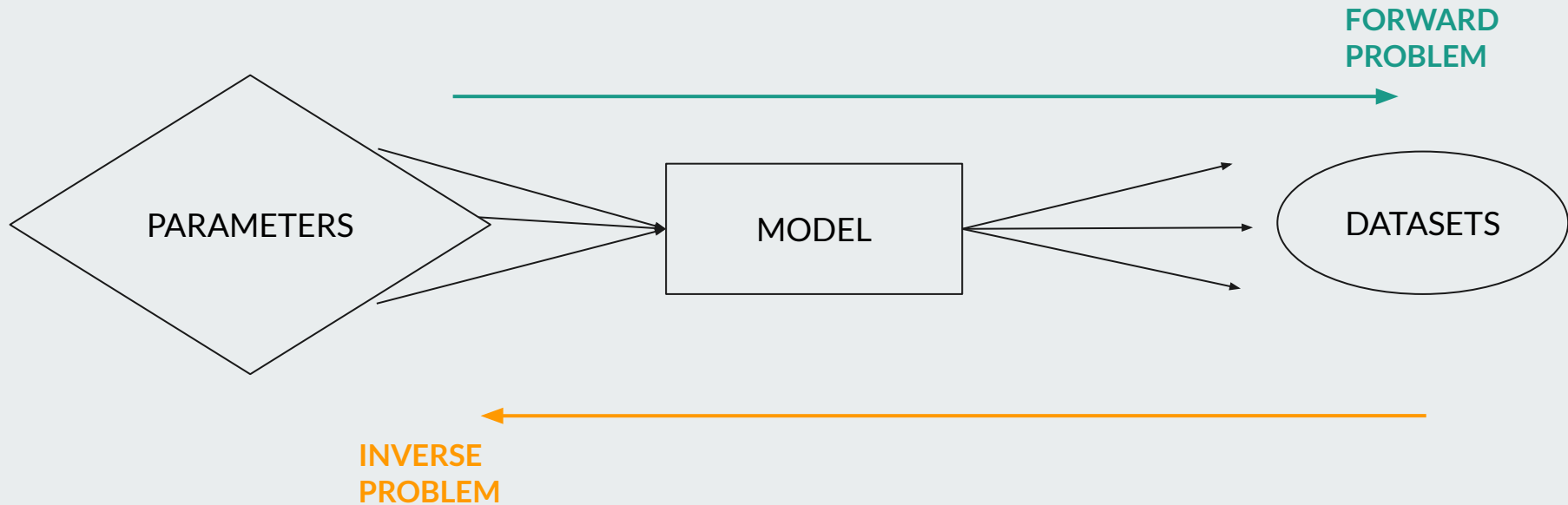




OK so far, you have a model and you can generate predictions given some parameters:



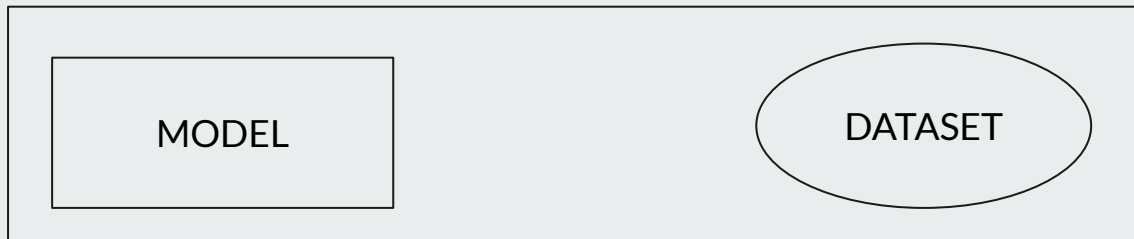


We can adopt a **qualitative** approach:

- Can my model reproduce some particular aspects of my dataset? (e.g. a positive relationship between two quantities, a hump-shaped trend for some variable...)
- For what range of parameters does this happen?

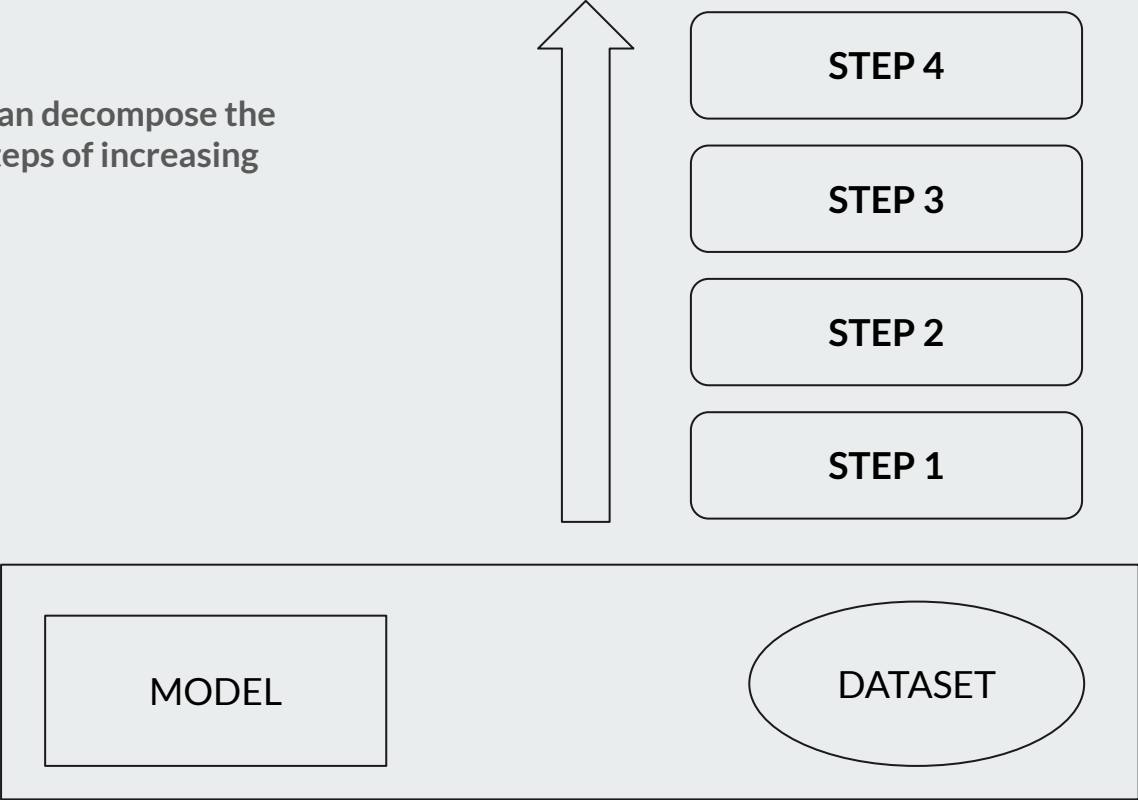
Generalized models can be useful here (e.g. Yeakel et al. 2011 *Theor. Ecol.*)

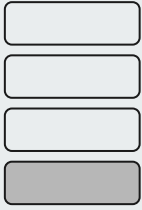
Or we we can adopt a more **quantitative** approach (stats)





For simplicity, we can decompose the process into four steps of increasing ambition:

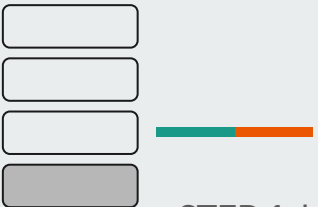




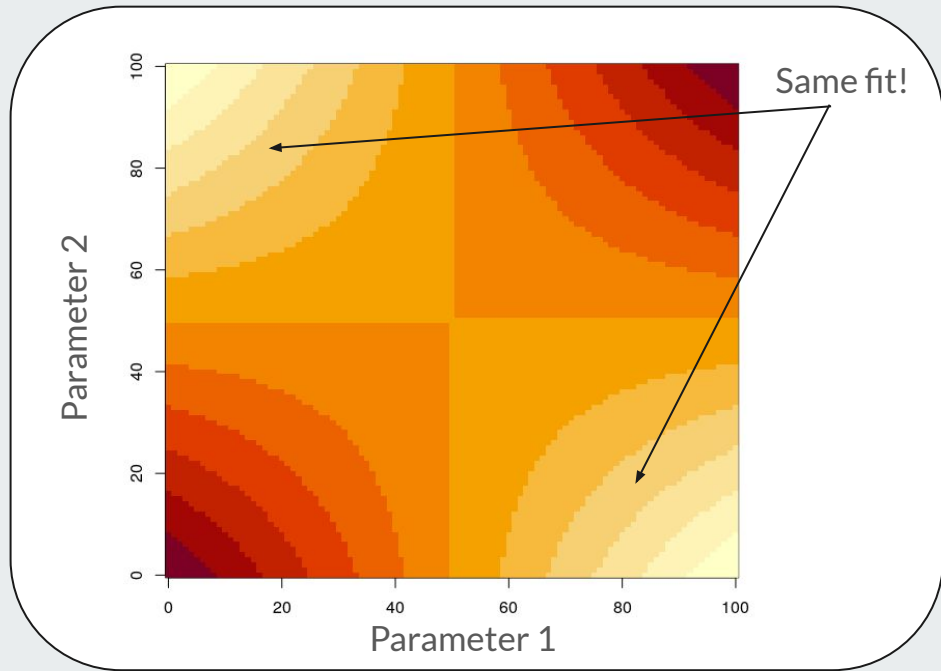
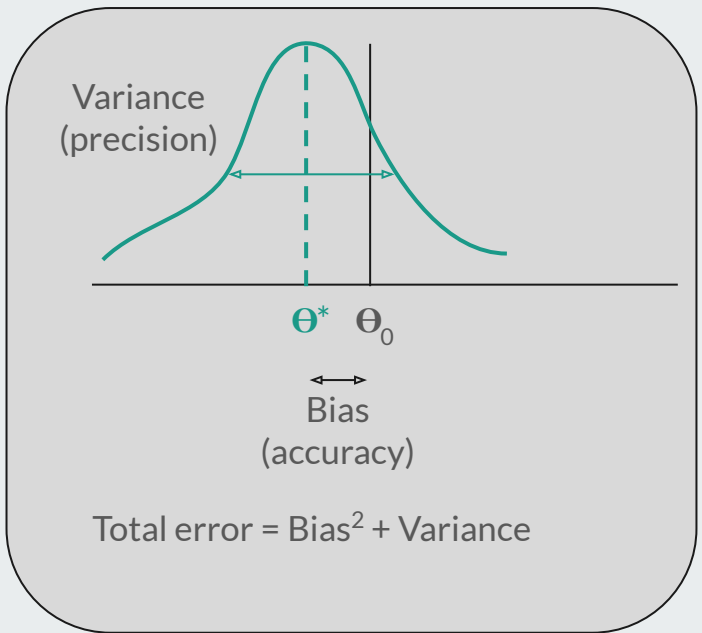
STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?

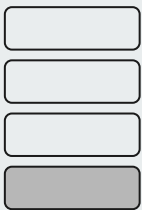
Major problems:

- we may obtain **biased** estimates
- the model may not be **identifiable**



STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?

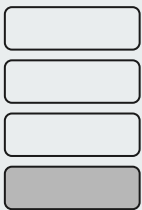




STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?

Common approaches:

- minimize some **distance** between predictions and data (e.g. least squares)
- maximize **likelihood**
- maximize **posterior probability** (Bayesian approaches)



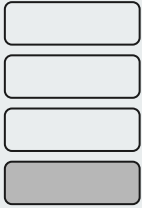
STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?

Least squares:

- + flexible, robust enough, very fast minimization techniques
- + equivalent to ML under certain assumptions*
- not applicable to all models
- somewhat ad-hoc: other distances could be used (absolute differences...)

Maximum likelihood:

- + fully general, intuitive, solid theoretical grounding
- + consistent (asymptotically unbiased... *if model is true*)
- can be hard to compute and maximize



STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?



IX. *On the Mathematical Foundations of Theoretical Statistics.*

By R. A. FISHER, M.A., *Fellow of Gonville and Caius College, Cambridge, Chief Statistician, Rothamsted Experimental Station, Harpenden.*

Communicated by DR. E. J. RUSSELL, F.R.S.

Received June 25,—Read November 17, 1921.



STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?

Least squares:

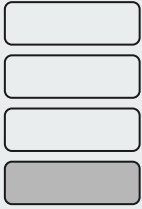
- + flexible, robust enough, very fast minimization techniques
- + equivalent to ML under certain assumptions*
- not applicable to all models
- somewhat ad-hoc: other distances could be used (absolute differences...)

Maximum likelihood:

- + fully general, intuitive, solid theoretical grounding
- + consistent (asymptotically unbiased... *if model is true*)
- can be hard to compute and maximize

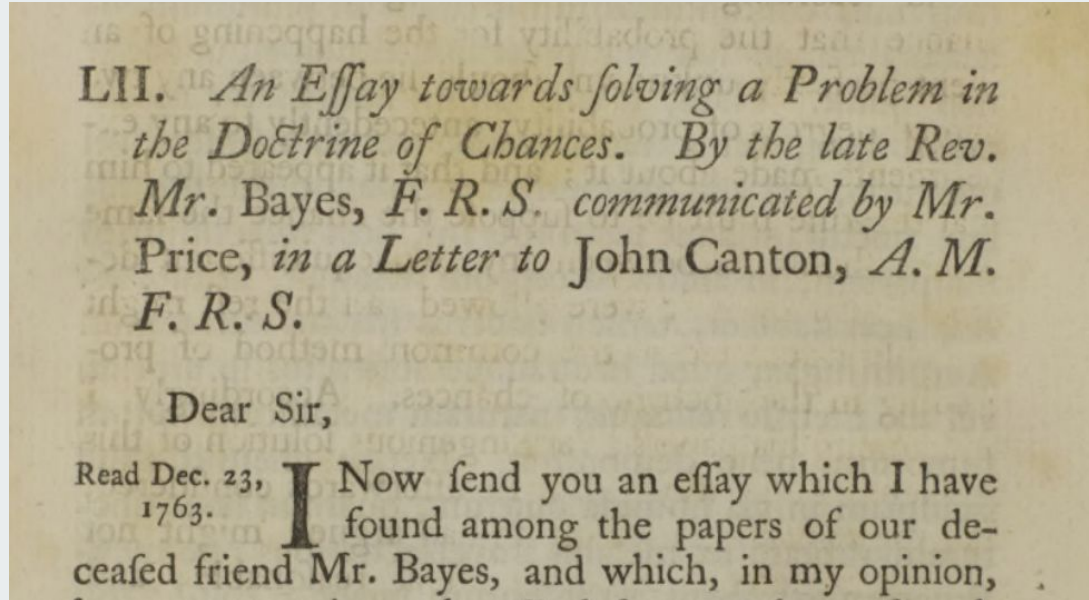
Posterior probability (Bayesian approaches):

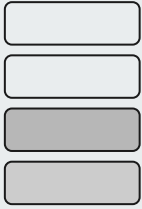
- + supplements ML with **prior knowledge** on parameter values
- + efficient sampling algorithms (priors guide the exploration of parameter space)
- + prior distributions can alleviate non-identifiability issues
- supplements ML with **prior knowledge** on parameter values
- can be slow to converge



STEP 1: how to find the best parameter values to fit the data (point estimation, model fitting)?

Bayes (1763)
Phil. Trans. Royal Soc.

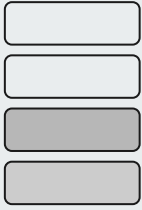




STEP 2: how to quantify the uncertainty in parameter estimates, and the quality of the fit (goodness of fit)?

Major problems:

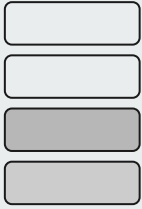
- we want to construct good **intervals** around parameter estimates: there are many of them
- with not so much data, we may have very little **statistical power** to evaluate the goodness of the fit (and failing to reject is not accepting)



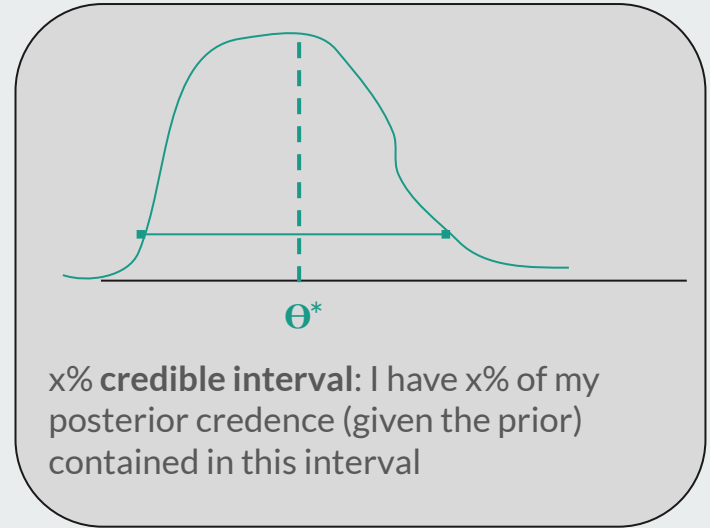
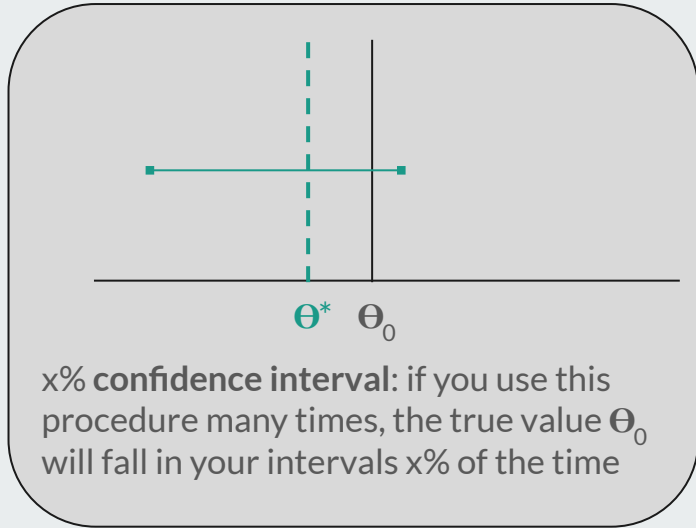
STEP 2: how to quantify the uncertainty in parameter estimates, and the quality of the fit (goodness of fit)?

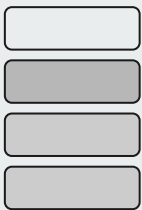
Common approaches (uncertainty):

- resample dataset/refit model (bootstrap, jackknife...)
- use likelihood surface theory to get confidence intervals (Fisher information...)
- use Bayesian approaches to compute credible intervals



STEP 2: how to quantify the uncertainty in parameter estimates, and the quality of the fit (goodness of fit)?

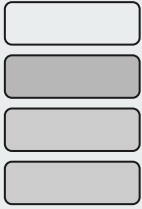




STEP 3: how to compare different models together and select the 'best' (model selection)?

Major problems:

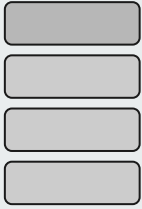
- a more complex model will always fit the data better but...
- the **bias/variance trade-off**, or the **curse of complexity**: for a given amount of data, too simple a model will have little variance/high bias (**underfitting**), too complex a model will have low bias/huge variance (**overfitting**)
- in both cases we have poor estimation of parameters (total uncertainty = $\text{bias}^2 + \text{variance}$)
- in both cases, we'll have poor prediction power for future/other datasets
- we must find some intermediate level of complexity, i.e. make some **compromise**



STEP 3: how to compare different models together and select the 'best' (model selection)?

Common approaches:

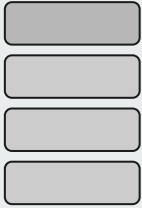
- Split dataset (**cross validation**, training/test datasets...)
- is a particular parameter 'significant' or not? (model simplification)
- **Likelihood ratio tests**
- Information criteria for model comparison (**AIC...**)
- Regularization or penalization techniques (ridge regression, LASSO...)



STEP 4: how to use several alternative models rather than just one (multimodel inference)?

Major problems:

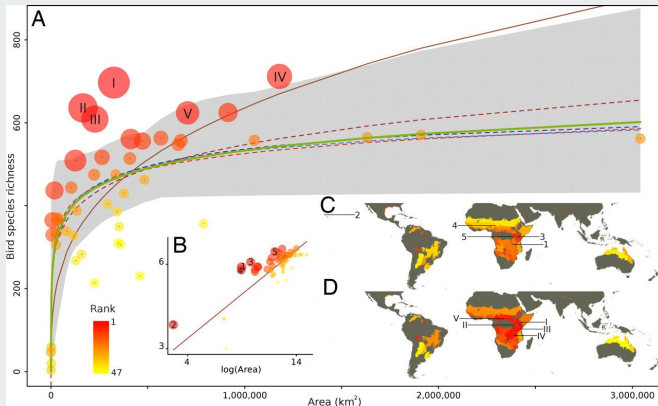
- How to reduce model selection bias (see Freedman's paradox)?
- How to include model selection uncertainty?
- How to combine the estimates or predictions from different competing models, and combine them in an optimal way?



STEP 4: how to use several alternative models rather than just one (multimodel inference)?

Common approaches:

- Take all models and do some ad-hoc **consensus** (e.g. average or median prediction)
- Use **model-averaging** techniques



Dormann et al. (2018) *Ecological Monographs*

SAR: Guilhaumon et al. (2008) *PNAS*
TPC: Padfield et al. (2021) *Methods in Ecology & Evolution*