# Introduction to Bayesian Statistics

# *Part 2* Bayesian Principles

Benjamin Rosenbaum



## **This lecture**

Frequentist statistics and Null hypothesis significance testing

Bayes' rule

Markov Chain Monte Carlo sampling

MCMC software

Why Bayesian statistics?

# **Statistical modeling**



# **Statistical modeling**

We are data detectives, trying to solve a mystery

In ecology, these mysteries can be extra tricky:

- Observational data instead of experiments
- Noisy data
- Many sources of variation

We're trying to unravel a signal from the noise (e.g. overall trend of biodiversity loss)

We want to make quantitative statements on research questions!



# Some comments on frequentist stats

# **Statistical modeling**



Maximum likelihood alone can't make probability statements about our research question!

# **Statistical modeling**



We want to assign probabilities to model parameters! E.g.,  $P(\theta > 0)$ 

- 1. We want to test hypothesis H1 (e.g., association is positive:  $\theta > 0$ )
- 2. We assume that the null hypothesis H0 is true (e.g., association is zero:  $\theta = 0$ )
- 3. Use a transformation T for which data Y (residuals) have a well-known distribution P(T) under H0
- 4. If T(Y) deviates enough from the assumed distribution, reject the null hypothesis  $P(T > T(Y) | \theta = 0) < 0.05 \rightarrow$  reject H0

Tests if the estimated association  $\theta^*$  (MLE) is just due to randomess of the data



- 1. We want to test hypothesis H1 (e.g., association is positive:  $\theta > 0$ )
- 2. We assume that the null hypothesis H0 is true (e.g., association is zero:  $\theta = 0$ )
- 3. Use a transformation T for which data Y (residuals) have a well-known distribution P(T) under H0
- 4. If T(Y) is improbable under assumed distribution, reject the null hypothesis  $P(T > T(Y) | \theta = 0) < 0.05 \rightarrow$  reject H0

Tests if the estimated association  $\theta^*$  (MLE) is just due to randomess of the data



- 1. We want to test hypothesis H1 (e.g., association is positive:  $\theta > 0$ )
- 2. We assume that the null hypothesis H0 is true
- (e.g., association is zero:  $\theta = 0$ ) **3. Use a transformation** *T* **for which data** *Y* **(residuals)** have a well-known distribution P(T) under H0
- 4. If T(Y) is improbable under assumed distribution, reject the null hypothesis  $P(T > T(Y) | \theta = 0) < 0.05 \rightarrow \text{reject H0}$

Tests if the estimated association  $\theta^*$  (MLE) is just due to randomess of the data



Transformation T

- 1. We want to test hypothesis H1 (e.g., association is positive:  $\theta > 0$ )
- 2. We assume that the null hypothesis H0 is true (e.g., association is zero:  $\theta = 0$ )
- 3. Use a transformation T for which data Y (residuals) have a well-known distribution P(T) under H0
- 4. If T(Y) is improbable under assumed distribution, reject the null hypothesis  $P(T > T(Y) | \theta = 0) < 0.05 \Rightarrow$  reject H0

Tests if the estimated association  $\theta^*$  (MLE) is just due to randomess of the data



Transformation T

# **Frequentist principles**

→ Data are a random realization of an experiment.
 True (but unknown) parameter is fixed.

Some problems with NHST:

- P-value: Probability of the data under the null hypothesis
- Can't confirm hypotheses, just reject the null hypothesis
- Standard errors rely on assumptions & approximations
- Confidence intervals' interpretation tricky
- Limited to tests with known distributions (T-test, F-test, ...)



50%-confidence intervals in 20 repeated experiments: 10 out of 20 contain the true value  $\mu$ 

# **Bayesian principles**

To make quantitative statements about research questions, we need probability distribution for model parameters, after observing the data  $P(\theta|y)$ 

 $\rightarrow$  Data is fixed, parameters are random.

Some examples:

- $P(\theta > 0) = 0.99$ : "I am 99% certain that the association is positive."
- 90%-quantile [0.5, 4.3]:
  "There is a 90% chance the slope is between 0.5 and 4.3."
- 2 population means  $\mu_1$  and  $\mu_2$ .  $P(\mu_1 - \mu_2)$  quantifies distribution of population-level difference.



Parameter θ

## **Bayesian models ???**

There is no such thing as a "Bayesian model"!

Maximum likelihood  $\rightarrow$  Point estimates

Frequentist NHST  $\rightarrow$  P-values for Null hypothesis

Bayesian stats

→ True probability distribution for model parameters

"Full luxury Bayes" (Richard McElreath)

"Bayesian 3D printer" (me)

#### The Bayesian 3D printer





Reverend Thomas Bayes (1701–1761)

Conditional probability: Prob. of event *A*, given that *B* occured

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

 $P(\text{Covid} \mid \text{positive}) = \frac{P(\text{positive} \mid \text{Covid}) \cdot P(\text{Covid})}{P(\text{positive})}$ 

P(positive | Covid)test sensitivityP(Covid)prevalence in the populationP(positive)positive test rate



Probability distribution of model parameters  $\theta$  after observing the data y



$$p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)}$$

#### $p(\theta|y)$ Posterior distribution

Update prior information in light of new evidence (data)

#### $p(\theta)$ Prior distribution

Belief about model parameters before data is observed

#### $p(y|\theta) = L(\theta)$ Likelihood function

Data inform the parameters, but L is not a probability distribution for  $\theta$ 

#### $p(y) = \int_{\theta} p(y|\theta) p(\theta)$ Normalization constant

Ensures that posterior is a true probility distribution with  $\int_{\theta} p(\theta|y) = 1$ 

Survival rate of a deer population 1 datapoint: 7 out of 10 individuals survived

Deterministic part:average survival rate  $\theta$ Stochastic part: $y \sim \text{Binomial}(N, \theta)$ 



Survival rate of a deer population 1 datapoint: 7 out of 10 individuals survived

Deterministic part:average survival rate  $\theta$ Stochastic part: $y \sim \text{Binomial}(N, \theta)$ 

**Prior:**  $\theta$  almost 0 or almost 1 are improbable  $\theta \sim \text{beta}(2,2)$  or  $p(\theta) = \text{dbeta}(\theta \mid 2,2)$ 



Survival rate of a deer population 1 datapoint: 7 out of 10 individuals survived

Deterministic part:average survival rate  $\theta$ Stochastic part: $y \sim \text{Binomial}(N, \theta)$ 

**Prior:**  $\theta$  almost 0 or almost 1 are improbable  $\theta \sim \text{beta}(2,2)$  or  $p(\theta) = \text{dbeta}(\theta \mid 2,2)$ 

**Likelihood:** defined by statistical model & data  $L(\theta) = p(y|\theta) = dBinomial(y = 7, N = 10 | \theta)$ 



Survival rate of a deer population 1 datapoint: 7 out of 10 individuals survived

Deterministic part:average survival rate  $\theta$ Stochastic part: $y \sim \text{Binomial}(N, \theta)$ 

**Prior:**  $\theta$  almost 0 or almost 1 are improbable  $\theta \sim \text{beta}(2,2)$  or  $p(\theta) = \text{dbeta}(\theta \mid 2,2)$ 

**Likelihood:** defined by statistical model & data  $L(\theta) = p(y|\theta) = dBinomial(y = 7, N = 10 | \theta)$ 

#### **Posterior:**

 $p(\theta|y) = \frac{L(\theta) \cdot p(\theta)}{c}$ 



## **Example: different prior**

Survival rate of a deer population 1 datapoint: 7 out of 10 individuals survived

Deterministic part:average survival rate  $\theta$ Stochastic part: $y \sim \text{Binomial}(N, \theta)$ 

**Prior:** uninformative  $\theta \sim beta(1,1) = uniform(1,1)$ 

**Likelihood:** defined by statistical model & data  $L(\theta) = p(y|\theta) = dBinomial(y = 7, N = 10 | \theta)$ 

#### Posterior:

 $p(\theta|y) = \frac{L(\theta) \cdot p(\theta)}{c}$  proportional to likelihood here



#### **Calculation of the posterior?**

To compute, e.g. mean( $\theta$ ), sd( $\theta$ ),  $P(\theta > 0.5)$  ...

we'd need to know  $p(\theta|y)$  for all  $\theta$  and also  $c = \int_{\theta} L(\theta) p(\theta)$ 

"Curse of

dimensionality"

1) Analytical (mathematical formula)

 $\rightarrow$  Much too complicated, often impossible

- 2) Numerical (e.g., grid)
- $\rightarrow$  Effort grows exponentially with  $\#\theta$
- $\rightarrow$  Computationally too expensive

**Oh no!** Same problem as before.



# Markov Chain Monte Carlo (MCMC) sampling

#### New idea: sampling!

$$p(\boldsymbol{\theta}|y) = \frac{p(y|\boldsymbol{\theta}) \cdot p(\boldsymbol{\theta})}{p(y)}$$

Instead of calculating  $p(\theta|y)$ , draw random  $\theta$  samples. Many samples where p high, few samples where p low

→ Sample density proportional to  $p(\theta|y)$ → We don't need the normalizing constant c = p(y)

 $p(\boldsymbol{\theta}|\boldsymbol{y}) \sim p(\boldsymbol{y}|\boldsymbol{\theta}) \cdot p(\boldsymbol{\theta})$ 

Posterior is proportional to likelihood x prior



#### **Markov Chain Monte Carlo**

Start with initial  $\theta_1$ Compute  $f(\theta_1) = L(\theta_1) \cdot p(\theta_1)$ 

In each step i = 2,3,4,...:

- Propose new  $\theta_{\text{new}}$ , e.g.  $\theta_{\text{new}} \sim \text{Normal}(\theta_{\text{old}}, \sigma)$ Compute  $f(\theta_{\text{new}}) = L(\theta_{\text{new}}) \cdot p(\theta_{\text{new}})$
- If  $f(\theta_{\text{new}}) > f(\theta_{\text{old}})$  $\rightarrow \text{accept } \theta_{i+1} = \theta_{\text{new}}$
- If  $f(\theta_{\text{new}}) < f(\theta_{\text{old}})$ 
  - → accept  $\theta_{i+1} = \theta_{\text{new}}$  with probability  $\frac{f(\theta_{\text{new}})}{f(\theta_{\text{old}})}$ (random draw)
  - $\rightarrow$  otherwise reject  $\theta_{\rm new}$



Repeat e.g. 1000 times

#### **Markov Chain Monte Carlo**

- $\theta_1, \theta_2, \theta_3, \dots, \theta_{1000}$  are called the **"chain**"
- They are samples from the underlying (but mathematically unknown) posterior distribution
- We can calculate empirical quantities
  - mean
  - standard deviation
  - quantiles
  - histogram (for visualization)
  - probability statements
- Without explicitly knowing the function  $p(\theta|y)$
- Even don't need to save computed values  $p(\theta_1|y), p(\theta_2|y), \dots$



Survival rate  $\theta$ 

## **Markov Chain Monte Carlo**

**"Markov"** property: each sample  $\theta_i$  **only** depends on the previous sample  $\theta_{i-1}$ 

**"Chain":** list of samples  $\theta_1, \theta_2, \theta_3, \dots, \theta_{1000}$ 

**"Monte Carlo":** each new sample involves a random draw

Very simple algorithm at its core (few lines of code)

Very sophisticated software to make it efficient (lots of maths go into good sample proposals)



Survival rate: 1 datapoint (7/10), prior beta(2,2)0.30 First sample  $\theta_1 = 0.4$ 0.20  $f(0.4) = dbinom(7,10, 0.4) \cdot dbeta(0.4, 2, 2)$ test = 0.0610.10 Propose  $\theta_{new} = 0.6$ 0.00  $f(0.6) = dbinom(7,10, 0.6) \cdot dbeta(0.6, 2, 2)$ = 0.310

 $f(\theta_{\text{new}}) > f(\theta_{\text{old}}) \rightarrow \text{accept } \theta_2 = 0.6 \text{ as next sample}$ 





Survival rate: 1 datapoint (7/10), prior beta(2,2)

Samples from posterior distribution  $\theta_1, \theta_2, \theta_3, ...$ describe probability for estimated survival rate. Don't need to know the full curve  $p(\theta|y)$  !

• Empirical **histogram** for visualization



Survival rate: 1 datapoint (7/10), prior beta(2,2)

Samples from posterior distribution  $\theta_1, \theta_2, \theta_3, ...$ describe probability for estimated survival rate. Don't need to know the full curve  $p(\theta|y)$  !

• Empirical mean and standard deviation

mean 
$$= \frac{1}{K} \sum_{i=1}^{K} \theta_i$$
  
sdev  $= \sqrt{\frac{1}{K} \sum_{i=1}^{K} (\theta_i - \text{mean})^2}$ 



Survival rate: 1 datapoint (7/10), prior beta(2,2)

Samples from posterior distribution  $\theta_1, \theta_2, \theta_3, ...$ describe probability for estimated survival rate. Don't need to know the full curve  $p(\theta|y)$  !

"Credible intervals"

90% of samples between the 5% and the 95% quantiles.

 $P(\theta \in [0.42, 0.81]) = 0.9$ 

",I am 90% sure the survival probability is between 0.42 and 0.81"



Survival rate  $\theta$ 

Survival rate: 1 datapoint (7/10), prior beta(2,2)

Samples from posterior distribution  $\theta_1, \theta_2, \theta_3, ...$ describe probability for estimated survival rate. Don't need to know the full curve  $p(\theta|y)$  !

• **Probability statements** (about hypotheses)

85% of samples larger that a survival rate of 0.5

 $P(\theta > 0.5) = 0.85$ 

"I am 85% sure the survival probability is larger than 0.5"



# **MCMC Demo**

#### https://chi-feng.github.io/mcmc-demo/app.html



- Mathematical theory says that MCMC will *eventually* be a good approximation of the posterior distribution
- How many samples are enough?
- Start with 1000-2000 samples
- Run multiple chains (3-4)
- Visual inspection
- Quantitative measures



Parameter  $\theta_1$ 

#### **Visual inspection**

- Traceplots for each parameter
- Should look like random noise
- Centered around a constant mean
- Chains should look similar
- Like a fuzzy caterpillar!

 $\rightarrow$  MCMC has converged



#### **Visual inspection**

- Traceplots for each parameter
- Should look like random noise
- Centered around a constant mean
- Chains should look similar
- Like a fuzzy caterpillar!

 $\rightarrow$  MCMC has converged



#### **Quantitative measures**

- Rhat value ("Gelman-Rubin statistic")
  - Compares the variation within and across chains
  - Value should be less than 1.1
- n\_eff (Number of effective samples)
  - Chains usually have a bit of autocorrelation but it shouldn't be too strong
  - Small n\_eff values indicate a problem

#### $\rightarrow$ MCMC has converged

```
3 chains with 1000 samples each
3000 samples (post-warmup)
Rhat = 1.001
n_eff= 1770
```





#### Some history

1700s Bayes' theorem, Laplace formalized it

Early 1800s Gauß: least squares, regression

**Late 1800s to early 1900s** Birth of modern statistics. Pearson, Fisher, Neyman ... : max. likelihood, hypothesis testing, design of experiments

Mid to late 1900s MCMC algorithms

**2000s** Computational tools for MCMC BUGS, JAGS, Stan ...

**Today** Convenient R interfaces brms, rstanarm ...

#### Future

Bayes impractical Restricted to simple cases

Frequentism superseded Bayes More practical in most cases

Still a niche topic in statistics

Becoming more popular in sciences

Taught in gradschools

Becoming the default instead of frequentism **??** 

#### Software

All Bayesian software contains:

#### 1) Modeling language

User must define statistical model:

parameters, likelihood & priors

#### 2) MCMC sampler

Automated algorithm that takes care of sampling

# **Bayesian programming languages**

- + Maximum flexibility in statistical modeling
- + Total control over every part of the model
- Steep learning curve
- Coding can be time-consuming

#### JAGS

Was the most popular once, now less and less used

Nimble

Extends JAGS, more flexible

Stan

Very efficient, runs in C++

Can all be called from **R** 

```
a ~ dnorm(0, sd=1) Nimble
b ~ dnorm(0, sd=1)
sigma ~ dexp(1.0)
for(i in 1:n) {
   y[i] ~ dnorm(a+b*x[i], sd=sigma)
}
```



## Formula-based R packages

- + Model formulation similar to Im or Ime4
- + Easy to learn
- + Less coding necessary
- + Handy functions for model analysis (after fitting)
- Limited to pre-defined model types
- "Im"-formulas deceive you into forgetting about model definition

# rstanarm GLMMs only Both automatically translate model into Stan Both automatically translate model into Stan brms Much more flexible, becoming quite popular Much more flexible, becoming quite popular Both automatically translate model into Stan brms Summary Stan Summary S

# **Specialized software** (R-packages)

- R-INLA, sdmTMB, → Spatial models spBayes
- bsam, hmmTMB, → Animal movement bayesmove
- sp0ccupancy → Occupancy models
- spAbundance  $\rightarrow$  Abundance models

...

• blavaan → Structural equation models



# Why Bayesian ?

# **Bayesian workflow**

- 1) Research question (hypotheses)
- 2) Data collection
- 3) Statistical model
- 4) Prior distribution choice
- 5) Model fitting (MCMC)
- 6) Evaluate model output
- 7) Quantitative statements on hypotheses



designed by 🕲 freepik.com

 $\rightarrow$  Workflow not that different from frequentist statistic

Revise model

# Why Bayesian?

#### **Philosophical answer:**

- Frequentism assumes true and fixed underlying parameter values.
- Data are just a sample of the "true" statistical model.
- Bayesian statistics embraces uncertainty and wants to quantify it correctly.
- Observed data are given, model parameters uncertain.



# Why Bayesian?

#### **Practical answer:**

- Output is more intuitive:
  - Direct inference on parameters / hypotheses instead of NHST What does the data tell me about my model?
- Full transparency and control over model and output
- Include prior belief / information
- Parameter regularization may be necessary
- Not limited to a specific toolbox, but full flexibility in modeling (especially with Stan or Nimble, but brms also very versatile)
- Fit complex models with lots of parameters



# Why Bayesian?

#### What are complex models?

- Nonlinear models
- Hierarchical structure, mixed effects
- Combination of multiple, heterogeneous data sources and/or models
- Constraints on parameters
- Latent variables

(occupancy models, animal movement models, SEM, HMM, ...)

#### The Bayesian 3D printer



#### Summary

- There is no such thing as a "Bayesian model"!
- Frequentist and Bayesian stats are different methodologies for estimating parameters of a statistical model
- **Frequentist** statistics cannot (mathematically) do direct inference  $P(\theta|y)$ , and requires a (methodological) detour via **NHST**  $P(y|\theta = 0)$
- **Bayesian** statistics can (conceptually) do direct inference  $P(\theta|y)$ , but requires a (computational) detour via **MCMC**

#### **Further reading**

Bürkner, P. (2024). The brms Book [in progress]. https://paulbuerkner.com/software/brms-book/

Fieberg, J. (2024). Statistics 4 Ecologists. <u>https://statistics4ecologists-v2.netlify.app/</u> [Chapters 11-13]

Inchausti, P. (2023). Statistical Modeling With R: a dual frequentist and Bayesian approach for life scientists. *Oxford University Press*. [Chapter 3]

Kery, M. & Kellner, F. (2024): Applied Statistical Modelling for Ecologists. *Elsevier*. [Chapter 2]

Johnson, A. A., Ott, M. Q., Dogucu, M. (2021). Bayes Rules! CRC Press. <u>https://www.bayesrulesbook.com/</u> [Chapters 1-2]

McElreath, R. (2020). Statistical Rethinking: A Bayesian Course with Examples in R and STAN (2nd ed.). *Chapman and Hall/CRC*. <u>https://doi.org/10.1201/9780429029608</u>

van de Schoot, R., Depaoli, S., King, R., et al. (2021). Bayesian statistics and modelling. *Nature Reviews. Methods Primers*, 1(1), 1–26. <u>https://doi.org/10.1038/s43586-020-00001-2</u>