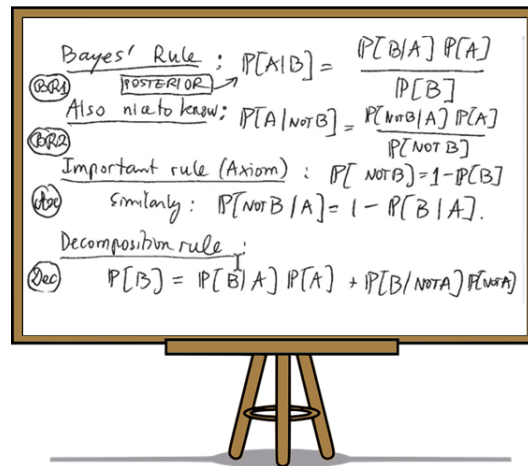


Bridging Theory and Practice Part 1: Introducing Bayesian inference with simple linear regression

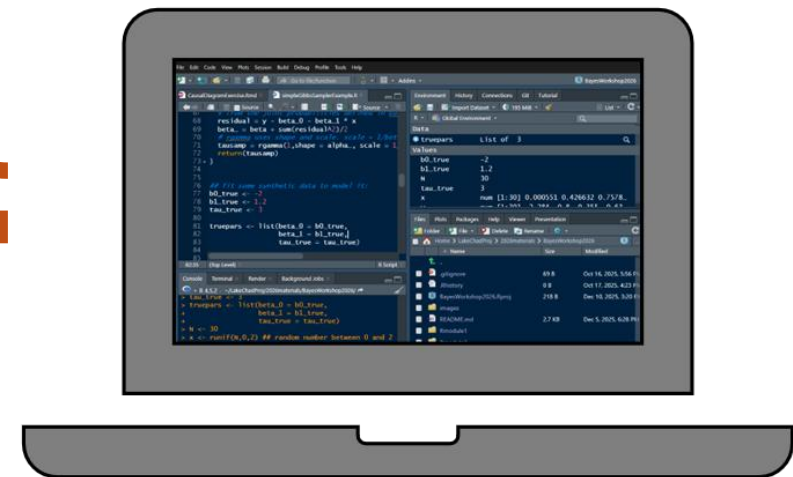
Katherine Muller



Theory



Bayes Rule!



Computation

Outline

- I. Origins of Bayesian inference
- II. Simple linear regression example
- III. Comparison with frequentist inference

Bayes' Theorem was first published in 1763

LII. *An Essay towards solving a Problem in the Doctrine of Chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S.*

Dear Sir,

Read Dec. 23, 1763. **I** Now send you an essay which I have found among the papers of our deceased friend Mr. Bayes, and which, in my opinion, has great merit, and well deserves to be preserved. Experimental philosophy, you will find, is nearly interested in the subject of it; and on this account there seems to be particular reason for thinking that a communication of it to the Royal Society cannot be improper.

“Bayesian inference” came later

Bayesian Analysis (2006)

1, Number 1, pp. 1–40

When Did Bayesian Inference Become “Bayesian”?

Stephen E. Fienberg*

Abstract. While Bayes' theorem has a 250-year history, and the method of inverse probability that flowed from it dominated statistical thinking into the twentieth century, the adjective “Bayesian” was not part of the statistical lexicon until relatively recently. This paper provides an overview of key Bayesian developments, beginning with Bayes' posthumously published 1763 paper and continuing up through approximately 1970, including the period of time when “Bayesian” emerged as the label of choice for those who advocated Bayesian methods.

Keywords: Bayes' Theorem; Classical statistical methods; Frequentist methods; Inverse probability; Neo-Bayesian revival; Stigler's Law of Eponymy; Subjective probability.

<https://projecteuclid.org/journals/bayesian-analysis/volume-1/issue-1/When-did-Bayesian-inference-become-Bayesian/10.1214/06-BA101.pdf>

What we call Bayes' Rule is elementary probability

a to N and that of b to N, i. e. the probability that the two subsequent events will both happen is compounded of the probability of the 1st and the probability of the 2d on supposition the 1st happens.

$$\longrightarrow P(A,B) = P(A) P(B|A)$$

By the same logic:

$$P(A,B) = P(B) P(A|B)$$

Therefore:

Rearrange:

$$\longleftarrow P(B) P(A|B) = P(A) P(B|A)$$

Bayes' Rule

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

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

Predict **unknown/unobservable** parameters based on **known/observable** data and prior understanding about parameters.

Fun fact: Bayes' friend Richard Price thought Bayes' theories confirmed the existence of God.

ry *. The purpose I mean is, to shew what reason we have for believing that there are in the constitution of things fixt laws according to which events happen, and that, therefore, the frame of the world must be the effect of the wisdom and power of an intelligent cause; and thus to confirm the argument taken from final causes for the existence of the Deity. It will be easy to see that the converse problem solved in this essay is more directly applicable to this purpose; for it shews us, with distinctness and precision, in every case of any particular order or recurrency of events, what reason there is to think that such recurrency or order is derived from stable causes or regulations innate, and not from any of the irregularities of chance.

your very humble servant,
Newington-Green,
Nov. 10, 1763.
Richard Price.
Ccc 2 SE C-

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

Predict **unknown/unobservable** parameters based on **known/observable** data and prior understanding about parameters.

Agroecology/sustainability science often involves estimating parameters from data

unknown/unobservable parameters

Disease resistance



Drought resilience



Effect of management
on soil health



known/observable data

Occurrence or severity of disease

Yields in dry and wet years

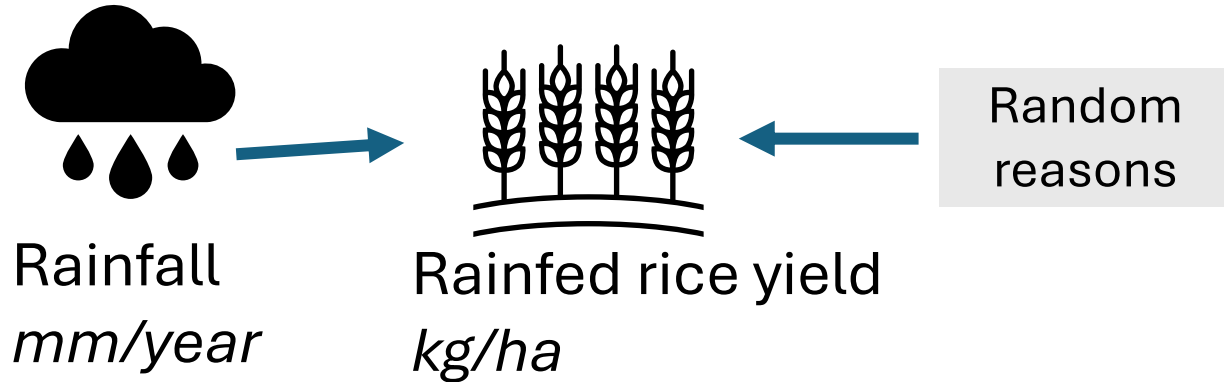
Management history and
soil health indicators

Estimating unknown parameters from data is not uniquely “Bayesian”

Part 2: Simple linear regression example

- Set up a model to explore theory and computation concepts
- Callback to causal inference and data story

Simple linear regression

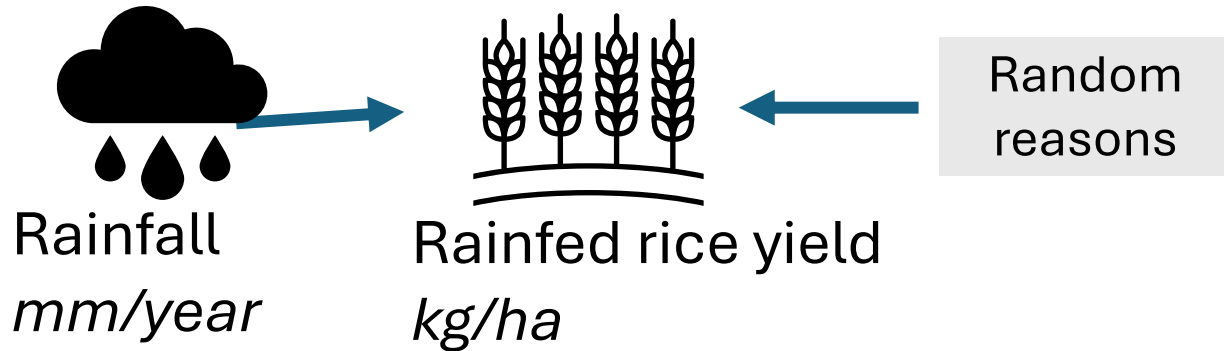


Yield = Function(Rainfall, Random)

Data story

- Rainfed rice is grown **on one farm** under optimal temperature, soil and management.
- Rice is planted and harvested **once per year every year** (no carryover).
- Rice yield is limited by **water availability through rainfall**.
- Rice yield varies randomly from year to year for reasons unrelated to rainfall

Simple linear regression



Yield = Function(Rainfall, Random)

Data processing: standardize variables

$$y = \frac{\text{yield} - \text{mean}(\text{yield})}{\text{sd}(\text{yield})}$$

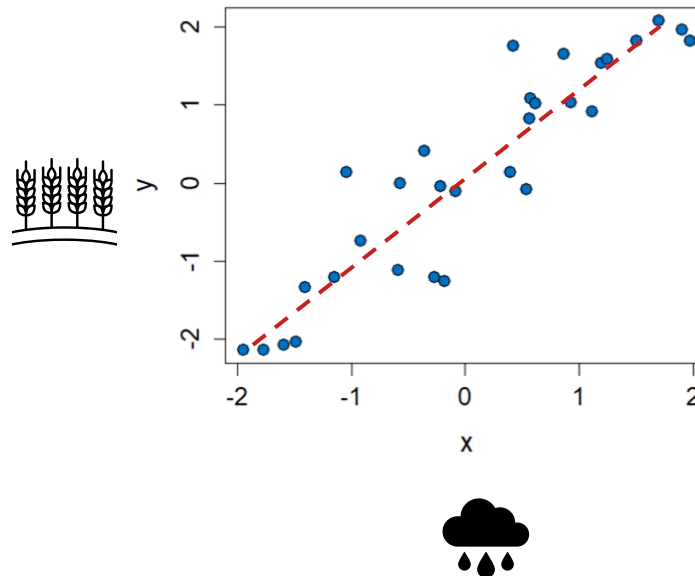
$$x = \frac{\text{rainfall} - \text{mean}(\text{rainfall})}{\text{sd}(\text{rainfall})}$$

- Two continuous random variables, **X**, and **Y**
- A change in **X** is associated with a change in **Y**.

$$y = \beta_0 + \beta_1 x + \varepsilon$$

intercept slope error

Error has Gaussian noise $\varepsilon \sim \mathcal{N}(0, 1/\tau)$



Precision (τ) is the reciprocal of variance

Data: x and y

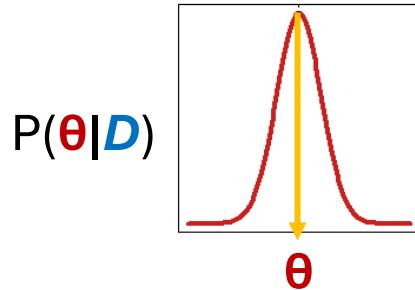
Parameters (θ): β_0, β_1, τ

Part 3: Frequentist vs. Bayesian inference

Overview

Frequentist inference

- Focus is on parameter **values**

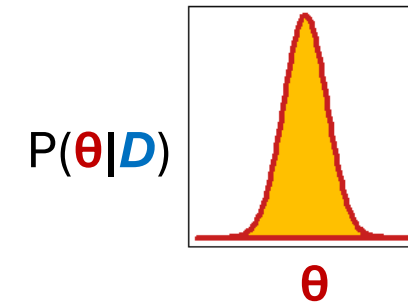


*With weak priors,
Bayesian and frequentist
methods give the same
parameter estimates.*

- Assumes **no prior understanding** about parameters
- Uncertainty** about parameters must be estimated indirectly, using resampling or asymptotic theory.

Bayesian inference

- Focus is on the parameter **distributions**



- Incorporates prior understanding** (or lack thereof).
- Uncertainty** about parameters is baked into the model

Frequentist Inference

Maximum Likelihood Estimation

Find **values** for each **parameter** that maximize the likelihood of the **data**

Likelihood function: $P(\mathbf{D} | \boldsymbol{\theta})$

Probability of the **data** given **parameters**

Easily solvable with math! (for simple models)

Simple linear regression:

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, 1/\tau)$$

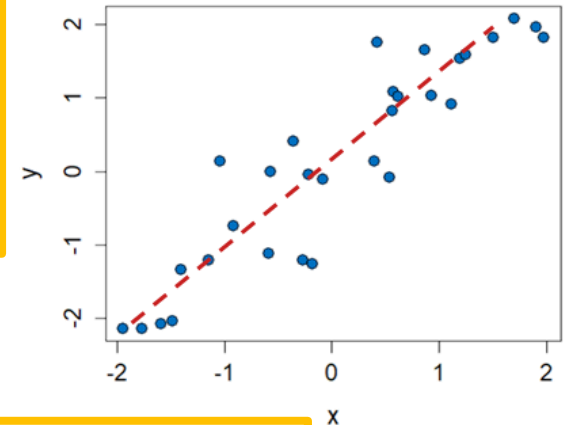
$$p(\mathbf{y}, \mathbf{x} | \beta_0, \beta_1, \tau) = \prod_{i=1}^n \mathcal{N}(\underbrace{\beta_0 + \beta_1 x_i}_{\text{line}}, \underbrace{1/\tau}_{\text{scatter}})$$

Find the values of β_0, β_1, τ that maximize this function with the observed \mathbf{x} and \mathbf{y} .

Equivalent for β_0 and β_1

Ordinary Least Squares:

Find the values of β_0, β_1 that draw the best fit line, minimizing the distance between line and points—i.e., sum of squared residuals

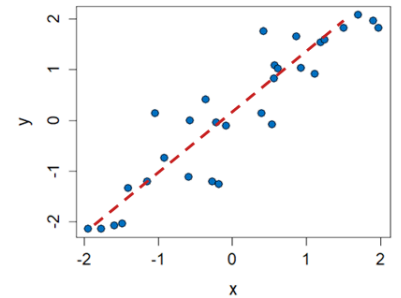


Bayesian Inference

Update our prior understanding about the **distribution** of **parameters** based on the **data**.

Simple linear regression:

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1/\tau)$$



$$\underbrace{P(\theta | D)}_{\text{Updated understanding about parameters (posterior)}} = \frac{\overbrace{P(D | \theta)}^{\text{Likelihood of the data given parameters}} \overbrace{P(\theta)}^{\text{Prior understanding about parameters}}}{\underbrace{P(D)}_{\text{Support over plausible values of the data: Normalizing constant so probability sums to 1}}}$$

$$p(\beta_0, \beta_1, \tau | \mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y} | \beta_0, \beta_1, \tau) p(\beta_0, \beta_1, \tau)}{p(\mathbf{x}, \mathbf{y})}$$

Find **most likely distributions** of β_0, β_1, τ based on observed \mathbf{x} and \mathbf{y} .

Bayesian Inference

Update our prior understanding about the **distribution** of **parameters** based on the **data**.

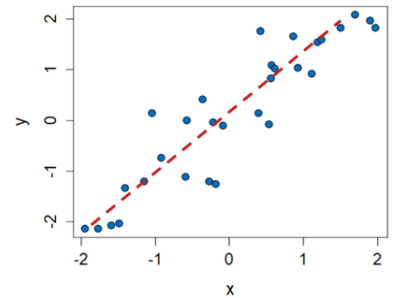
Simple linear regression:

$$y = \beta_0 + \beta_1 x + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1/\tau)$$

Specify prior distributions for parameters

$$P(\theta) \quad \begin{array}{l} \beta_0 \sim \\ \beta_1 \sim \\ \tau \sim \end{array} \boxed{f(\theta^*)} \quad \begin{array}{l} \text{intercept} \\ \text{slope} \\ \text{precision} \end{array}$$

hyperparameters



Choose priors based on:

- Analytical convenience
- How they represent your knowledge (or lack of knowledge) about your system

$$p(\beta_0, \beta_1, \tau \mid x, y) = \frac{p(x, y \mid \beta_0, \beta_1, \tau) p(\beta_0, \beta_1, \tau)}{p(x, y)}$$

Find **most likely distributions** of β_0, β_1, τ based on observed x and y .

Bayesian Inference

Update our prior understanding about the **distribution** of **parameters** based on the **data**.

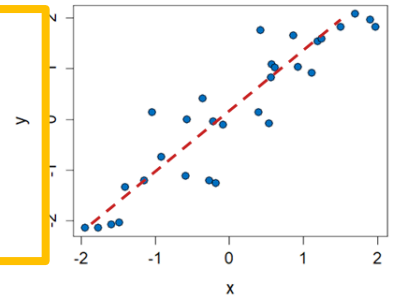
Simple linear regression:

$$y = \beta_0 + \beta_1 x + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1/\tau)$$

Specify prior distributions for parameters

$$P(\theta) \quad \begin{aligned} \beta_0 &\sim \mathcal{N}(\mu_0, (v_0 \tau)^{-1}) \text{ intercept} \\ \beta_1 &\sim \mathcal{N}(\mu_1, (v_1 \tau)^{-1}) \text{ slope} \\ \tau &\sim \text{Gamma}(a, b) \text{ precision} \end{aligned}$$

Prior uncertainty about the regression coefficients must be proportional to prediction uncertainty for the data



Choose priors based on:

- Analytical convenience

$$p(\beta_0, \beta_1, \tau | x, y) = \frac{p(x, y | \beta_0, \beta_1, \tau) p(\beta_0, \beta_1, \tau)}{p(x, y)}$$

Find **most likely distributions** of β_0, β_1, τ based on observed x and y .

With very specific priors, the posterior can be solved analytically (without computer sampling)

Bayesian Inference

Update our prior understanding about the **distribution** of **parameters** based on the **data**.

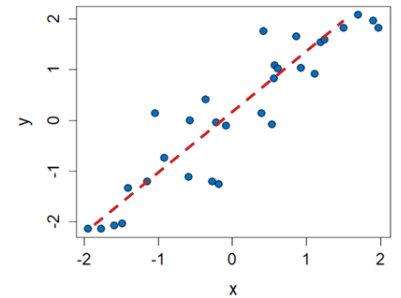
Simple linear regression:

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, 1/\tau)$$

Specify prior distributions for parameters

$$\begin{aligned} P(\boldsymbol{\theta}) \quad & \beta_0 \sim \mathcal{N}(\mu_0, \tau_0^{-1}) && \text{intercept} \\ & \beta_1 \sim \mathcal{N}(\mu_1, \tau_1^{-1}) && \text{slope} \\ & \tau \sim \text{Gamma}(\mathbf{a}, \mathbf{b}) && \text{precision} \end{aligned}$$

Independent hyperparameters



$$p(\beta_0, \beta_1, \tau \mid \mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y} \mid \beta_0, \beta_1, \tau) p(\beta_0, \beta_1, \tau)}{p(\mathbf{x}, \mathbf{y})}$$

Find **most likely distributions** of β_0, β_1, τ based on observed \mathbf{x} and \mathbf{y} .

The posterior can no longer be solved without computer sampling

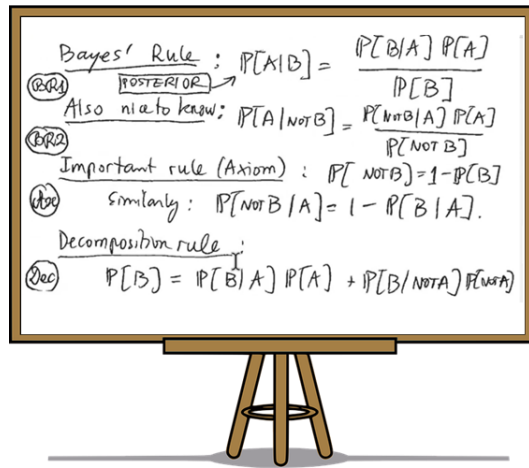
Recap

- Statistical inference (Bayesian and otherwise) involves estimating **unknown or unobservable parameters** from **observable data** (e.g., **disease resistance** from **disease occurrence**).
- Frequentist inference differs from Bayesian in that it:
 1. Focuses on parameter values instead of distributions
 2. Estimates parameters solely from the data (no priors)
 3. Is less well-suited to dealing with uncertainty about parameters
 4. Is typically easier and faster to compute.
- Frequentist and Bayesian methods produce similar estimates for parameter values when prior expectations are weak.
- Without computer sampling, Bayesian inference can only be performed on simple models with very restricted priors.

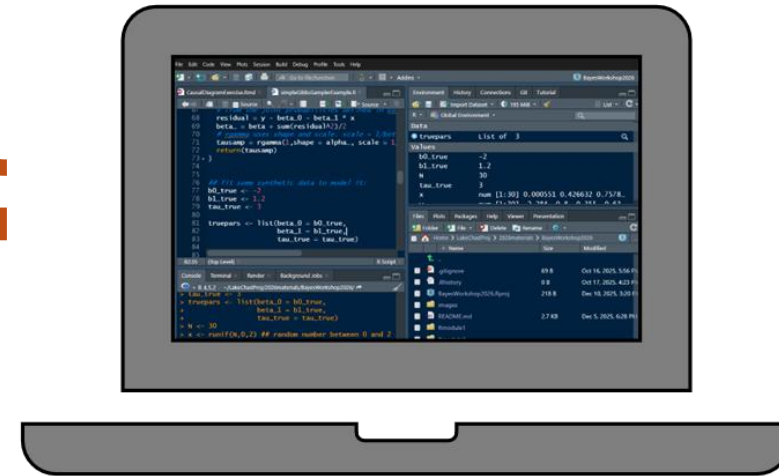
Frequentist vs. Bayesian inference: Computation

Aspect	Frequentist MLE	Bayesian
Goal	Find the best parameters that fit the data	Estimate the full distribution of parameters
Computation	One optimization run	Many repeated calculations (sampling or approximation)
Speed / Effort	Usually fast, moderate computation	Can be slow, especially for complex models
Determinism	Gives the same result each time	Results can vary due to randomness in sampling
Ease	Easy to do with standard tools	Harder; often requires specialized software

Next: Why computation makes Bayesian inference tractable



Theory



Computation