

Bayesian Modeling and Inference: An Introduction to STAN for the Social Sciences

Workshop for Social Scientists

Dr Diego Perez Ruiz

University of Manchester
School of Social Statistics

May 23, 2025

Course Outline for Today

Morning Session (HBS G7)

- 09:30 – 10:00: Introduction and Session Aims
- 10:00 – 11:30: Theoretical – Bayes (Part 1)
 - What is Bayesian Statistics? Examples.
 - Bayes' Theorem: From Events to Parameters
 - Priors, Likelihood, Posterior
 - Example: Normal Data, Unknown Mean
 - Point & Interval Estimation (Credible Intervals)
- 11:30 – 11:50: Break
- 11:50 – 12:10: Theoretical – Bayes (Part 2)

Lunch Break (12:30 – 1:30)

Afternoon Session (HBS 2.88) -

Practical

- 1:30 – 4:15: Practical Application
 - Tools (Stan overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison (LOO, WAIC, DIC)
- 4:15 – 4:30: Review and Closing

Outline for this Section

- 1 Introduction and Session Aims (9:30 – 10:00)
- 2 Theoretical – Bayes (Part 1) (10:00 – 11:30)
 - What is Bayesian Statistics?
 - Statistical Distributions
 - Priors
 - Example: Normal Data, Unknown Mean (Known Variance)
- 3 Theoretical – Bayes (Part 2) (11:50am – 12:10pm)
 - The Computational Challenge Intro to MCMC
- 4 Summary of Morning Session (12:10 – 12:30pm)
- 5 Practical Application (1:30 – 4:15pm)
 - Tools (Stan Overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison
- 6 Review and Closing (4:15 – 4:30pm)

Bayesian Modeling and Inference: An Introduction to STAN for the Social Sciences

Dr Diego Perez Ruiz
School of Social Statistics
University of Manchester

May 23, 2025

- Understand the fundamental concepts of Bayesian statistics.
- Differentiate Bayesian approaches from frequentist statistics.
- Learn about key components: prior distributions, likelihood, posterior distributions.
- Get introduced to the Bayesian workflow for model building and evaluation.
- Gain a conceptual understanding of Markov Chain Monte Carlo (MCMC) and its diagnostics.
- Be introduced to Stan as a tool for Bayesian inference.
- See practical examples of linear and logistic regression using Stan.
- Understand methods for model comparison (LOO-CV, WAIC, DIC).

- All practical materials are on GitHub:
 - github.com/dapr12/BayesSTAN-SocialSci

- All practical materials are on GitHub:
 - github.com/dapr12/BayesSTAN-SocialSci
- The repository includes a README.md with full setup instructions.

- All practical materials are on GitHub:
 - github.com/dapr12/BayesSTAN-SocialSci
- The repository includes a README.md with full setup instructions.
- Each tutorial section below details its specific files:
 - A **PDF guide** for instructions theory.
 - An **R script** for code execution.
 - Necessary **Stan model files** ('.stan').

Contents: Tutorial 1 - Bayesian Linear Regression

- **PDF Guide:** `bayesian_linear_regression_tutorial.pdf`
 - Step-by-step instructions, explanations, and exploration questions.

Contents: Tutorial 1 - Bayesian Linear Regression

- **PDF Guide:** `bayesian_linear_regression_tutorial.pdf`
 - Step-by-step instructions, explanations, and exploration questions.
- **R Script:** `bayesian_linear_regression.R`
 - Contains all R code from the PDF for you to run.

Contents: Tutorial 1 - Bayesian Linear Regression

- **PDF Guide:** `bayesian_linear_regression_tutorial.pdf`
 - Step-by-step instructions, explanations, and exploration questions.
- **R Script:** `bayesian_linear_regression.R`
 - Contains all R code from the PDF for you to run.
- **Stan Model:** `linear_regression.stan`
 - Defines the Bayesian linear regression model structure.

Contents: Tutorial 1 - Bayesian Linear Regression

- **PDF Guide:** `bayesian_linear_regression_tutorial.pdf`
 - Step-by-step instructions, explanations, and exploration questions.
- **R Script:** `bayesian_linear_regression.R`
 - Contains all R code from the PDF for you to run.
- **Stan Model:** `linear_regression.stan`
 - Defines the Bayesian linear regression model structure.
- **Focus:** Data simulation, model fitting, MCMC diagnostics, PPCs, and interpretation for linear models.

Contents: Tutorial 2 - Bayesian Logistic Regression

- **PDF Guide:** `bayesian_logistic_regression_tutorial.pdf`
 - Comprehensive guide for logistic regression (binary outcomes, logit link).

Contents: Tutorial 2 - Bayesian Logistic Regression

- **PDF Guide:** `bayesian_logistic_regression_tutorial.pdf`
 - Comprehensive guide for logistic regression (binary outcomes, logit link).
- **R Script:** `bayesian_logistic_regression.R`
 - All R code examples from the logistic regression PDF.

Contents: Tutorial 2 - Bayesian Logistic Regression

- **PDF Guide:** `bayesian_logistic_regression_tutorial.pdf`
 - Comprehensive guide for logistic regression (binary outcomes, logit link).
- **R Script:** `bayesian_logistic_regression.R`
 - All R code examples from the logistic regression PDF.
- **Stan Model:** `logistic_regression.stan`
 - Defines the Bayesian logistic regression model structure.

Contents: Tutorial 2 - Bayesian Logistic Regression

- **PDF Guide:** `bayesian_logistic_regression_tutorial.pdf`
 - Comprehensive guide for logistic regression (binary outcomes, logit link).
- **R Script:** `bayesian_logistic_regression.R`
 - All R code examples from the logistic regression PDF.
- **Stan Model:** `logistic_regression.stan`
 - Defines the Bayesian logistic regression model structure.
- **Focus:** Modeling binary data and log-odds/odds ratio interpretation.

Contents: Tutorial 3 - Bayesian Model Comparison

- **PDF Guide:** `bayesian_model_comparison_tutorial.pdf`
 - Explains model comparison using PSIS-LOO and WAIC with examples.

Contents: Tutorial 3 - Bayesian Model Comparison

- **PDF Guide:** `bayesian_model_comparison_tutorial.pdf`
 - Explains model comparison using PSIS-LOO and WAIC with examples.
- **R Script:** `bayesian_model_comparison.R`
 - R code to fit multiple models and perform comparisons using the 'loo' package.

Contents: Tutorial 3 - Bayesian Model Comparison

- **PDF Guide:** `bayesian_model_comparison_tutorial.pdf`
 - Explains model comparison using PSIS-LOO and WAIC with examples.
- **R Script:** `bayesian_model_comparison.R`
 - R code to fit multiple models and perform comparisons using the 'loo' package.
- **Stan Models:**
 - `linear_model_x1.stan` (predicts y from x_1)
 - `linear_model_x1_x2.stan` (predicts y from x_1 and x_2)

Contents: Tutorial 3 - Bayesian Model Comparison

- **PDF Guide:** `bayesian_model_comparison_tutorial.pdf`
 - Explains model comparison using PSIS-LOO and WAIC with examples.
- **R Script:** `bayesian_model_comparison.R`
 - R code to fit multiple models and perform comparisons using the 'loo' package.
- **Stan Models:**
 - `linear_model_x1.stan` (predicts y from x_1)
 - `linear_model_x1_x2.stan` (predicts y from x_1 and x_2)
- **Focus:** Calculating ELPD, PSIS-LOO, WAIC, effective parameters.

Why Bayesian Statistics for Social Sciences?

Why Bayesian Statistics for Social Sciences?

Why Bayesian Statistics for Social Sciences?

- **Intuitive Interpretation:** Probability as a degree of belief.

Why Bayesian Statistics for Social Sciences?

- **Intuitive Interpretation:** Probability as a degree of belief.
- **Incorporating Prior Knowledge:** Formally use existing research or theory.

Why Bayesian Statistics for Social Sciences?

- **Intuitive Interpretation:** Probability as a degree of belief.
- **Incorporating Prior Knowledge:** Formally use existing research or theory.
- **Rich Uncertainty Quantification:** Full distributions, not just p-values.

Why Bayesian Statistics for Social Sciences?

- **Intuitive Interpretation:** Probability as a degree of belief.
- **Incorporating Prior Knowledge:** Formally use existing research or theory.
- **Rich Uncertainty Quantification:** Full distributions, not just p-values.
- **Model Flexibility:** Build models that match complex social theories (hierarchical, etc.).

Why Bayesian Statistics for Social Sciences?

- **Intuitive Interpretation:** Probability as a degree of belief.
- **Incorporating Prior Knowledge:** Formally use existing research or theory.
- **Rich Uncertainty Quantification:** Full distributions, not just p-values.
- **Model Flexibility:** Build models that match complex social theories (hierarchical, etc.).
- **Small Sample Sizes:** Priors can help when data is limited.

Why Bayesian Statistics for Social Sciences?

- **Intuitive Interpretation:** Probability as a degree of belief.
- **Incorporating Prior Knowledge:** Formally use existing research or theory.
- **Rich Uncertainty Quantification:** Full distributions, not just p-values.
- **Model Flexibility:** Build models that match complex social theories (hierarchical, etc.).
- **Small Sample Sizes:** Priors can help when data is limited.
- MCMC methods are generally used on Bayesian models for fitting realistically complex models.

A man and his tools make a man and his trade — Vita Sackville-West

We shape our tools and then the tools shape us — Winston Churchill

Tools: Probabilistic Programming Languages (PPLs)

- **Stan** (BSD-3 License) - *Our primary focus*

Why PPLs?

These tools automate the complex MCMC sampling process, allowing us to focus on model specification and interpretation.

Tools: Probabilistic Programming Languages (PPLs)

- **Stan** (BSD-3 License) - *Our primary focus*
- **Turing.jl** (MIT License) - Julia-based PPL

Why PPLs?

These tools automate the complex MCMC sampling process, allowing us to focus on model specification and interpretation.

Tools: Probabilistic Programming Languages (PPLs)

- **Stan** (BSD-3 License) - *Our primary focus*
- **Turing.jl** (MIT License) - Julia-based PPL
- PyMC (Apache License) - Python-based

Why PPLs?

These tools automate the complex MCMC sampling process, allowing us to focus on model specification and interpretation.

Tools: Probabilistic Programming Languages (PPLs)

- **Stan** (BSD-3 License) - *Our primary focus*
- **Turing.jl** (MIT License) - Julia-based PPL
- PyMC (Apache License) - Python-based
- JAGS (GPL License) - "Just Another Gibbs Sampler"

Why PPLs?

These tools automate the complex MCMC sampling process, allowing us to focus on model specification and interpretation.

Tools: Probabilistic Programming Languages (PPLs)

- **Stan** (BSD-3 License) - *Our primary focus*
- **Turing.jl** (MIT License) - Julia-based PPL
- PyMC (Apache License) - Python-based
- JAGS (GPL License) - "Just Another Gibbs Sampler"
- BUGS (GPL License) - "Bayesian inference Using Gibbs Sampling" (an early PPL)

Why PPLs?

These tools automate the complex MCMC sampling process, allowing us to focus on model specification and interpretation.



- **High-performance** platform for statistical modeling and computation.



- **High-performance** platform for statistical modeling and computation.
- Strong financial support from **NUMFocus** and industry (AWS, Bloomberg, Microsoft, IBM, RStudio, Facebook, NVIDIA, Netflix).



- **High-performance** platform for statistical modeling and computation.
- Strong financial support from **NUMFocus** and industry (AWS, Bloomberg, Microsoft, IBM, RStudio, Facebook, NVIDIA, Netflix).
- Open-source language, syntax **similar to C++**.



- **High-performance** platform for statistical modeling and computation.
- Strong financial support from **NUMFocus** and industry (AWS, Bloomberg, Microsoft, IBM, RStudio, Facebook, NVIDIA, Netflix).
- Open-source language, syntax **similar to C++**.
- Uses state-of-the-art MCMC algorithms (primarily **NUTS - No-U-Turn Sampler**).



- **High-performance** platform for statistical modeling and computation.
- Strong financial support from **NUMFocus** and industry (AWS, Bloomberg, Microsoft, IBM, RStudio, Facebook, NVIDIA, Netflix).
- Open-source language, syntax **similar to C++**.
- Uses state-of-the-art MCMC algorithms (primarily **NUTS - No-U-Turn Sampler**).
- Offers interfaces from **R (rstan)**, Python (CmdStanPy, PyStan), Julia, etc.



- **High-performance** platform for statistical modeling and computation.
- Strong financial support from **NUMFocus** and industry (AWS, Bloomberg, Microsoft, IBM, RStudio, Facebook, NVIDIA, Netflix).
- Open-source language, syntax **similar to C++**.
- Uses state-of-the-art MCMC algorithms (primarily **NUTS - No-U-Turn Sampler**).
- Offers interfaces from **R (rstan)**, Python (CmdStanPy, PyStan), Julia, etc.
- MCMC **parallel sampler** (runs multiple chains in parallel for efficiency).

Stan Code Example: Simple Linear Regression

```
1 data { // Data provided to Stan
2   int<lower=0> N;      // Number of observations
3   vector[N] x;        // Predictor variable
4   vector[N] y;        // Outcome variable
5 }
6 parameters { // Parameters to be estimated
7   real alpha;         // Intercept
8   real beta;          // Slope
9   real<lower=0> sigma; // Error standard deviation (must be
    positive)
10 }
11 model { // Priors and Likelihood
12   // Priors
13   alpha ~ normal(0, 20);
14   beta ~ normal(0, 2);
15   sigma ~ cauchy(0, 2.5); // Half-Cauchy for sigma
16
17   // Likelihood
18   y ~ normal(alpha + beta * x, sigma);
19 }
20
```

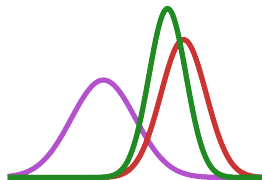
Listing 1: Conceptual Stan model structure.

Stan Code Example: Simple Linear Regression

```
1 data { // Data provided to Stan
2   int<lower=0> N;      // Number of observations
3   vector[N] x;        // Predictor variable
4   vector[N] y;        // Outcome variable
5 }
6 parameters { // Parameters to be estimated
7   real alpha;         // Intercept
8   real beta;          // Slope
9   real<lower=0> sigma; // Error standard deviation (must be
    positive)
10 }
11 model { // Priors and Likelihood
12   // Priors
13   alpha ~ normal(0, 20);
14   beta ~ normal(0, 2);
15   sigma ~ cauchy(0, 2.5); // Half-Cauchy for sigma
16
17   // Likelihood
18   y ~ normal(alpha + beta * x, sigma);
19 }
20
```

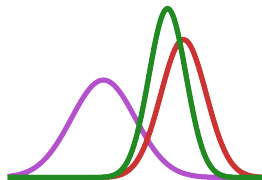
Listing 2: Conceptual Stan model structure.

- An ecosystem of Julia packages for Bayesian inference.



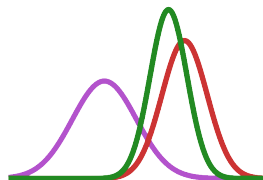
Alternative PPL: Turing.jl

- An **ecosystem of Julia packages** for Bayesian inference.
- Written in **Julia**, compiled with LLVM (often leading to high performance).



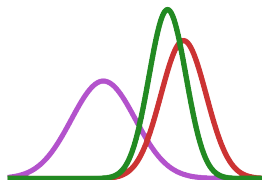
Alternative PPL: Turing.jl

- An **ecosystem of Julia packages** for Bayesian inference.
- Written in **Julia**, compiled with LLVM (often leading to high performance).
- Financial support from **NUMFocus**.



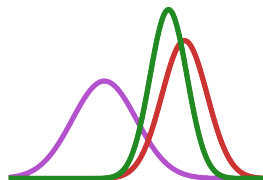
Alternative PPL: Turing.jl

- An **ecosystem of Julia packages** for Bayesian inference.
- Written in **Julia**, compiled with LLVM (often leading to high performance).
- Financial support from **NUMFocus**.
- Highly **composable** with other Julia packages.



Alternative PPL: Turing.jl

- An **ecosystem of Julia packages** for Bayesian inference.
- Written in **Julia**, compiled with LLVM (often leading to high performance).
- Financial support from **NUMFocus**.
- Highly **composable** with other Julia packages.
- Offers **several MCMC algorithms** (NUTS, HMC, Particle Gibbs, etc.).



Turing.jl Ecosystem Highlights

Key packages within the Turing ecosystem include:

- `Turing.jl`: The main interface for defining models.

Composability

The modular nature of the Julia ecosystem allows these tools to work together seamlessly.

Turing.jl Ecosystem Highlights

Key packages within the Turing ecosystem include:

- `Turing.jl`: The main interface for defining models.
- `MCMCChains.jl`: For diagnostics, summaries, and plotting of MCMC output.

Composability

The modular nature of the Julia ecosystem allows these tools to work together seamlessly.

Turing.jl Ecosystem Highlights

Key packages within the Turing ecosystem include:

- `Turing.jl`: The main interface for defining models.
- `MCMCChains.jl`: For diagnostics, summaries, and plotting of MCMC output.
- `DynamicPPL.jl`: Provides the domain-specific language (DSL) for probabilistic programming.

Composability

The modular nature of the Julia ecosystem allows these tools to work together seamlessly.

Turing.jl Ecosystem Highlights

Key packages within the Turing ecosystem include:

- `Turing.jl`: The main interface for defining models.
- `MCMCChains.jl`: For diagnostics, summaries, and plotting of MCMC output.
- `DynamicPPL.jl`: Provides the domain-specific language (DSL) for probabilistic programming.
- `AdvancedHMC.jl`: Advanced Hamiltonian Monte Carlo algorithms.

Composability

The modular nature of the Julia ecosystem allows these tools to work together seamlessly.

Turing.jl Ecosystem Highlights

Key packages within the Turing ecosystem include:

- `Turing.jl`: The main interface for defining models.
- `MCMCChains.jl`: For diagnostics, summaries, and plotting of MCMC output.
- `DynamicPPL.jl`: Provides the domain-specific language (DSL) for probabilistic programming.
- `AdvancedHMC.jl`: Advanced Hamiltonian Monte Carlo algorithms.
- `DistributionsAD.jl`: Automatic differentiation for log-PDFs of distributions.

Composability

The modular nature of the Julia ecosystem allows these tools to work together seamlessly.

Turing.jl Ecosystem Highlights

Key packages within the Turing ecosystem include:

- `Turing.jl`: The main interface for defining models.
- `MCMCChains.jl`: For diagnostics, summaries, and plotting of MCMC output.
- `DynamicPPL.jl`: Provides the domain-specific language (DSL) for probabilistic programming.
- `AdvancedHMC.jl`: Advanced Hamiltonian Monte Carlo algorithms.
- `DistributionsAD.jl`: Automatic differentiation for log-PDFs of distributions.
- `Bijectors.jl`: Handles transformations for constrained variables.

Composability

The modular nature of the Julia ecosystem allows these tools to work together seamlessly.

Turing.jl Code Example: Simple Linear Regression

```
1 using Turing, Distributions
2
3 @model function linreg(x, y)
4     # Priors
5     alpha ~ Normal(0, 20)
6     beta ~ Normal(0, 2)
7     # For positive sigma, use truncated or a distribution defined on
       R+
8     sigma ~ truncated(Cauchy(0, 2.5); lower=0)
9
10    # Likelihood (vectorized)
11    # The '.' before '~' and '.*' indicates broadcasting
12    y .~ Normal(alpha .+ beta .* x, sigma)
13 end
14
15 # Example usage (conceptual):
16 # N = length(x_data)
17 # model_instance = linreg(x_data, y_data)
18 # chain = sample(model_instance, NUTS(), 1000)
19
```

Listing 3: Conceptual Turing.jl model structure.

Turing.jl Code Example: Simple Linear Regression

- Models are defined within a '@model' macro as Julia functions.

Turing.jl Code Example: Simple Linear Regression

- Models are defined within a '@model' macro as Julia functions.
- Syntax is often more concise and closer to statistical notation.

Outline for this Section

- 1 Introduction and Session Aims (9:30 – 10:00)
- 2 Theoretical – Bayes (Part 1) (10:00 – 11:30)
 - What is Bayesian Statistics?
 - Statistical Distributions
 - Priors
 - Example: Normal Data, Unknown Mean (Known Variance)
- 3 Theoretical – Bayes (Part 2) (11:50am – 12:10pm)
 - The Computational Challenge Intro to MCMC
- 4 Summary of Morning Session (12:10 – 12:30pm)
- 5 Practical Application (1:30 – 4:15pm)
 - Tools (Stan Overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison
- 6 Review and Closing (4:15 – 4:30pm)

Recommended References

- Gelman et al. (2013) - Chapter 1: Probability and inference
- McElreath (2020) - Chapter 1: The Golem of Prague
- Gelman, Hill and Vehtari (2020) - Chapter 3: Basic methods in math probability
- Khan and Rue (2021)
- **Probability Textbooks:**
 - Bertsekas and Tsitsiklis (2008)
 - Dekking et al. (2010) (skip frequentist parts)
 - Jaynes (2003) (philosophical)
 - Kurt (2019) (playful)
 - Diaconis and Skyrms (2019) (philosophical, less rigor)

A Little Motivation

Inside every nonBayesian there is a Bayesian struggling to get out
— Denis Lindley

What is Bayesian Statistics?

What is Bayesian Statistics?

What is Bayesian Statistics?

Definition

Bayesian statistics is a **data analysis approach based on Bayes' theorem** where available knowledge about the parameters of a statistical model is updated with the information of observed data. (Gelman et al., 2013).

What is Bayesian Statistics?

Definition

Bayesian statistics is a **data analysis approach based on Bayes' theorem** where available knowledge about the parameters of a statistical model is updated with the information of observed data. (Gelman et al., 2013).

- Observable quantities x (the data) vs. Unknown quantities θ (parameters).
- Parameters are treated as **random variables**.
- We make probability statements about model parameters.

Bayes' Theorem for Events

Named after Rev. Thomas Bayes (1702-1761).

For probability events A and B

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

Bayes' Theorem for Events

Named after Rev. Thomas Bayes (1702-1761).

For probability events A and B

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

For a set of mutually exclusive and exhaustive events A_i

$$p(A_i|B) = \frac{p(B|A_i)p(A_i)}{\sum_j p(B|A_j)p(A_j)}$$

Example: Coin Tossing

Let A = event of 2 Heads in three tosses of a fair coin. Let B = event the 1st coin toss is a Head.

- Sample space: $\{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$
(8 outcomes)
- $A = \{HHT, HTH, THH\} \implies p(A) = 3/8$
- $B = \{HHH, HHT, HTH, HTT\} \implies p(B) = 1/2$
- $p(B|A) = p(A \cap B)/p(A) = (2/8)/(3/8) = 2/3$

Example: Coin Tossing

Let A = event of 2 Heads in three tosses of a fair coin. Let B = event the 1st coin toss is a Head.

- Sample space: $\{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$ (8 outcomes)
- $A = \{HHT, HTH, THH\} \implies p(A) = 3/8$
- $B = \{HHH, HHT, HTH, HTT\} \implies p(B) = 1/2$
- $p(B|A) = p(A \cap B)/p(A) = (2/8)/(3/8) = 2/3$

Using Bayes' Theorem:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)} = \frac{(2/3) \cdot (3/8)}{(1/2)} = 1/2$$

Example: Diagnostic Testing (HIV)

A new HIV test is claimed to have “95% sensitivity and 98% specificity”. In a population with an HIV prevalence of 1/1000, what is the chance that a patient testing positive actually has HIV?

Example: Diagnostic Testing (HIV)

- A : Patient is truly HIV positive (has the disease).
 - $P(A) = 0.001$ (This is the disease prevalence, our **prior** probability).
- A' : Patient is truly HIV negative (does not have the disease).
 - $P(A') = 1 - P(A) = 1 - 0.001 = 0.999$.
- B : Patient's test result is positive.

Example: Diagnostic Testing (HIV)

- A : Patient is truly HIV positive (has the disease).
 - $P(A) = 0.001$ (This is the disease prevalence, our **prior** probability).
- A' : Patient is truly HIV negative (does not have the disease).
 - $P(A') = 1 - P(A) = 1 - 0.001 = 0.999$.
- B : Patient's test result is positive.
- **Sensitivity:** $P(B|A) = 0.95$
 - "If a person **has** HIV, the test correctly identifies them as positive 95% of the time." (True Positive Rate)

Example: Diagnostic Testing (HIV)

- A : Patient is truly HIV positive (has the disease).
 - $P(A) = 0.001$ (This is the disease prevalence, our **prior** probability).
- A' : Patient is truly HIV negative (does not have the disease).
 - $P(A') = 1 - P(A) = 1 - 0.001 = 0.999$.
- B : Patient's test result is positive.
- **Sensitivity:** $P(B|A) = 0.95$
 - "If a person **has** HIV, the test correctly identifies them as positive 95% of the time." (True Positive Rate)
- **Specificity is 98%**, which means:
 - "If a person **does not have** HIV, the test correctly identifies them as negative 98% of the time." (True Negative Rate)
 - Therefore, the **False Positive Rate** is $P(B|A') = 1 - \text{Specificity} = 1 - 0.98 = 0.02$.
 - "If a person **does not have** HIV, there's a 2% chance the test will incorrectly say they are positive."

Example: Diagnostic Testing (HIV)

- A : Patient is truly HIV positive (has the disease).
 - $P(A) = 0.001$ (This is the disease prevalence, our **prior** probability).
- A' : Patient is truly HIV negative (does not have the disease).
 - $P(A') = 1 - P(A) = 1 - 0.001 = 0.999$.
- B : Patient's test result is positive.
- **Sensitivity:** $P(B|A) = 0.95$
 - "If a person **has** HIV, the test correctly identifies them as positive 95% of the time." (True Positive Rate)
- **Specificity is 98%**, which means:
 - "If a person **does not have** HIV, the test correctly identifies them as negative 98% of the time." (True Negative Rate)
 - Therefore, the **False Positive Rate** is $P(B|A') = 1 - \text{Specificity} = 1 - 0.98 = 0.02$.
 - "If a person **does not have** HIV, there's a 2% chance the test will incorrectly say they are positive."

Example: Diagnostic Testing (HIV)

We want to find $P(A|B)$: the probability the patient has HIV given they tested positive.

$$\begin{aligned}P(A|B) &= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A')P(A')} \\&= \frac{0.95 \times 0.001}{(0.95 \times 0.001) + (0.02 \times 0.999)} \\&= \frac{0.00095}{0.00095 + 0.01998} = \frac{0.00095}{0.02093} \approx 0.045389... \\&\approx 4.5\% \quad (\text{This is our } \mathbf{posterior} \text{ probability})\end{aligned}$$

Example: Diagnostic Testing (HIV)

We want to find $P(A|B)$: the probability the patient has HIV given they tested positive.

$$\begin{aligned} P(A|B) &= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A')P(A')} \\ &= \frac{0.95 \times 0.001}{(0.95 \times 0.001) + (0.02 \times 0.999)} \\ &= \frac{0.00095}{0.00095 + 0.01998} = \frac{0.00095}{0.02093} \approx 0.045389... \\ &\approx 4.5\% \quad (\text{This is our } \mathbf{posterior} \text{ probability}) \end{aligned}$$

Surprising Result!

Even with a positive test, there's only about a 4.5% chance the patient actually has HIV in this low-prevalence population! Our prior belief (low prevalence) strongly influences the posterior.

Being Bayesian: From Events to Parameters

- The HIV example shows how a test result (data) *updates our prior belief*.
- Bayesian statistical inference applies this to unknown *parameters* θ .

Bayes' Theorem for Parameters

$$\underbrace{p(\theta|x)}_{\text{Posterior}} = \frac{\overbrace{p(x|\theta)}^{\text{Likelihood}} \cdot \overbrace{p(\theta)}^{\text{Prior}}}{\underbrace{\int p(x|\theta')p(\theta')d\theta'}_{\text{Marginal Likelihood / Evidence}}}$$

Often written: **Posterior** \propto **Likelihood** \times **Prior**

- $p(\theta)$: **Prior distribution** - uncertainty about θ *before* data.
- $p(x|\theta)$: **Likelihood** - model of data given θ .
- $p(\theta|x)$: **Posterior distribution** - uncertainty about θ *after* data.

What changes from Frequentist Statistics?

- **Flexibility** - probabilistic building blocks (priors, likelihoods)¹.
- **Better uncertainty treatment**: Full posterior distributions.
- **Intuitive Interpretation**: "95% probability θ is in $[a,b]$ " (Credible Interval).
- No p-values (in the traditional sense).

¹like LEGO

Frequentist Definition of Probability

- Based on observation of a large number of trials (long-run frequency).
- For an event E , if we get n_E successes out of n trials, then the probability $P(E)$ is:

$$P(E) = \lim_{n \rightarrow \infty} \frac{n_E}{n}$$

- Probability is seen as an objective property of the phenomenon, measurable through repeated experiments.

Bayesian Definition of Probability

- Probability $P(E)$ reflects our *degree of belief* or *state of knowledge* about event E .
- This belief is expressed as a probability distribution.
- This distribution must be consistent with all of our existing beliefs.
- **Example of Inconsistency (Not Allowed):**
 - Believing $P(\text{Coin is Heads}) = 0.7$
 - AND Believing $P(\text{Coin is Tails}) = 0.8$ (if only Heads/Tails possible and these must sum to 1)
- Probability is subjective (relative to the observer's knowledge) but must be coherent.

Frequentist Approach: Confidence Level & Intervals

- A confidence level describes the *procedure* of constructing intervals.
- It's the proportion of random samples (from the same population) that would produce confidence intervals containing the true, fixed (but unknown) population parameter.
- **Example:** A 95% confidence level means:
 - If we were to draw 100 random samples from the population,
 - And construct a 95% confidence interval from each sample,
 - Approximately 95 of these 100 intervals would contain the true population parameter.

Frequentist Confidence Interval: Crucial Point

- For any *single, specific* confidence interval calculated from one sample:
 - The true parameter either IS in that interval or IS NOT.
 - The probability that this **particular** interval contains the true parameter is either 0 or 1.
- **Therefore:** The confidence level (e.g., 95%) is NOT the probability that a given, already calculated interval includes the true population parameter.

Example 1.6: Frequentist Confidence Interval

Poll Context (2015 Pew Research, 1,500 adults):

"We are 95% confident that 60% to 64% of Americans think the federal government does not do enough for middle class people."

Example 1.6: Frequentist Confidence Interval

Poll Context (2015 Pew Research, 1,500 adults):

"We are 95% confident that 60% to 64% of Americans think the federal government does not do enough for middle class people."

The Correct Frequentist Interpretation:

- "If we were to repeat this polling process many times, 95% of the random samples of 1,500 adults would produce confidence intervals (like [60%, 64%]) that contain the true proportion of Americans who hold this belief."

Misconception 1:

- *“There is a 95% chance that this specific confidence interval [60%, 64%] includes the true population proportion.”*

Frequentist CI: Common Misconceptions for [60%, 64%]

Misconception 1:

- *"There is a 95% chance that this specific confidence interval [60%, 64%] includes the true population proportion."*
- **Why incorrect (Frequentist view):** This interval is fixed based on our one sample. The true proportion is fixed. This interval either contains it (Prob=1) or it doesn't (Prob=0). We just don't know which.

Misconception 2:

Frequentist CI: Common Misconceptions for [60%, 64%]

Misconception 1:

- *"There is a 95% chance that this specific confidence interval [60%, 64%] includes the true population proportion."*
- **Why incorrect (Frequentist view):** This interval is fixed based on our one sample. The true proportion is fixed. This interval either contains it (Prob=1) or it doesn't (Prob=0). We just don't know which.

Misconception 2:

- *"The true population proportion is in this interval [60%, 64%] 95% of the time."*

Frequentist CI: Common Misconceptions for [60%, 64%]

Misconception 1:

- *“There is a 95% chance that this specific confidence interval [60%, 64%] includes the true population proportion.”*
- **Why incorrect (Frequentist view):** This interval is fixed based on our one sample. The true proportion is fixed. This interval either contains it (Prob=1) or it doesn't (Prob=0). We just don't know which.

Misconception 2:

- *“The true population proportion is in this interval [60%, 64%] 95% of the time.”*
- **Why incorrect:** The true population proportion is a fixed, constant value. It doesn't move in and out of intervals. It's the intervals (constructed from different hypothetical samples) that would vary.

Frequentist View on a Single CI

- To a frequentist, for a specific, calculated confidence interval, the probability that it captures the true parameter is either zero or one.
- The challenge is that one *never knows* whether that specific interval is one of the ones that contains the true value (probability one) or one that doesn't (probability zero).
- So, a frequentist makes a statement about the *long-run performance of the method used to construct intervals*:
 "95% of similarly constructed intervals contain the true value."

The Bayesian Alternative: Credible Interval

- Bayesian credible intervals have a definition that is often considered more intuitive and easier to interpret directly.
- Since a Bayesian is allowed to express uncertainty about parameters in terms of probability (using a posterior distribution):
 - A Bayesian credible interval is a range derived from the posterior distribution.
 - For this range, the Bayesian believes the probability of it including the true parameter value is, for example, 0.95 (or 95%).
- **Key Difference:** A Bayesian *can* make a direct probabilistic statement about a specific interval containing the true parameter value.

Bayesian Credible Interval: Direct Interpretation

Continuing the example idea:

If a Bayesian analysis yielded a 95% credible interval of [60%, 64%] for the proportion of Americans who think the federal government does not do enough for middle class people:

Bayesian Credible Interval: Direct Interpretation

Continuing the example idea:

If a Bayesian analysis yielded a 95% credible interval of [60%, 64%] for the proportion of Americans who think the federal government does not do enough for middle class people:

- **A Bayesian can state:**

“Given our model and the observed data, there is a 95% probability that the true proportion of Americans with this belief lies between 60% and 64%.”

- This interpretation aligns more closely with what many people intuitively **want** a confidence interval to mean.

On the Nature of Probability



Bruno de Finetti

Yes, probability does not exist (as an objective physical quantity)...

- If we disregard objective chance, nothing is lost.
- The math of inductive rationality remains the same.

Subjective Probability

Probability is a degree of belief, given a state of information.

Recommended References

- Grimmett and Stirzaker (2020): Ch 3 (Discrete), Ch 4 (Continuous)
- Dekking et al. (2010): Ch 4 (Discrete), Ch 5 (Continuous)
- Betancourt (2019)

Probability Distributions: The Building Blocks

- Bayesian statistics uses probability distributions as the inference engine.
- Think of them as "Lego" pieces to construct models.

Probability Distribution Function (PDF/PMF)

- Mathematical function outputting probabilities (discrete) or densities (continuous).
- $P(X) : X \rightarrow \mathbb{R}$ (integral / sum is 1).
- Discrete: Probability Mass Function (PMF).
- Continuous: Probability Density Function (PDF).

$$X \sim \text{Dist}(\theta_1, \theta_2, \dots)$$

- X : random variable.
- Dist : distribution name (e.g., Normal, Bernoulli).
- $\theta_1, \theta_2, \dots$: parameters (e.g., mean, std. dev.).

Example

$$Y \sim \text{Normal}(\mu, \sigma)$$

Priors

Recommended References

- Gelman et al. (2013): Ch 2 (Single-param), Ch 3 (Multi-param)
- McElreath (2020) - Ch 4: Geocentric Models
- Gelman, Hill Vehtari (2020): Ch 9 (Prior info, types of priors)
- Schoot et al. (2021)

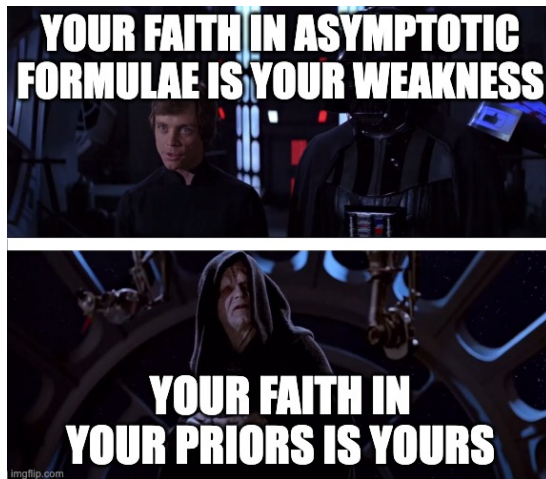


Figure: The power of the prior!

Prior Probability

The prior $P(\theta)$ in Bayes' Theorem:

$$P(\theta|y) \propto P(y|\theta) \cdot P(\theta)$$

Represents beliefs about parameters *before* seeing the data.

The Subjectivity of the Prior

- A common criticism.
- **However:** All modeling involves subjective choices.
- Bayesian stats makes these choices **explicit and formal** via priors.
- Frequentist choices (model form, error structure) are also subjective but often less transparent.

Types of Priors

- **Uniform (flat):** Often not recommended. Can be improper.
 - Example: $N(0, \text{huge variance})$ like $N(0, 10000)$, approximates $U(-\infty, \infty)$.
- **Weakly informative:** Gentle regularization. Good default.
- **Informative:** Based on strong external knowledge. Use with justification.

Typically, we specify that we have no prior knowledge ... We therefore specify vague, diffuse or uninformative priors.

Weakly Uninformative Prior Examples

- For regression coefficients (standardized predictors):
 - $\theta \sim \text{Normal}(0, s)$ (e.g., $s = 1, 2.5$)
 - $\theta \sim \text{Student-t}(\nu = 3, 0, s = 2.5)$ (robust)

Weakly Uninformative Prior Examples

- For regression coefficients (standardized predictors):
 - $\theta \sim \text{Normal}(0, s)$ (e.g., $s = 1, 2.5$)
 - $\theta \sim \text{Student-t}(\nu = 3, 0, s = 2.5)$ (robust)
- For variance parameters σ (must be positive):
 - Half-Cauchy, Half-Normal, Half-Student-t, Exponential.

The "Normal" Trap

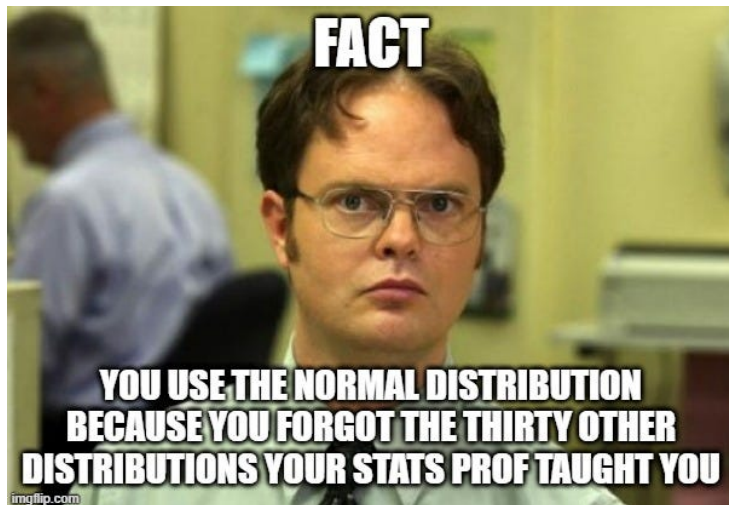


Figure: Don't forget other distributions!

Examples of Bayesian Inference using the Normal Distribution

Focus of this Section

We will now explore concrete examples of how Bayesian inference works, specifically when dealing with data assumed to come from a Normal distribution.

The Normal distribution is foundational in statistics, and understanding Bayesian approaches with it provides a strong basis for more complex models.

Known Variance, Unknown Mean: Model Setup

It is easier to first consider a model with only one unknown parameter.

Scenario:

- Suppose we have a sample of Normal data: $x_i \sim N(\mu, \sigma^2)$, for $i = 1, \dots, n$.
- We assume the variance, σ^2 , is **known**.
- The mean, μ , is **unknown**.

Known Variance, Unknown Mean: Model Setup

It is easier to first consider a model with only one unknown parameter.

Scenario:

- Suppose we have a sample of Normal data: $x_i \sim N(\mu, \sigma^2)$, for $i = 1, \dots, n$.
- We assume the variance, σ^2 , is **known**.
- The mean, μ , is **unknown**.

Prior Distribution for the Mean μ :

- We assume a prior distribution for μ based on our prior beliefs.
- Let this prior be Normal: $\mu \sim N(\mu_0, \sigma_0^2)$.
 - μ_0 is the prior mean.
 - σ_0^2 is the prior variance, reflecting our initial uncertainty about μ .

Known Variance, Unknown Mean: Model Setup

It is easier to first consider a model with only one unknown parameter.

Scenario:

- Suppose we have a sample of Normal data: $x_i \sim N(\mu, \sigma^2)$, for $i = 1, \dots, n$.
- We assume the variance, σ^2 , is **known**.
- The mean, μ , is **unknown**.

Prior Distribution for the Mean μ :

- We assume a prior distribution for μ based on our prior beliefs.
- Let this prior be Normal: $\mu \sim N(\mu_0, \sigma_0^2)$.
 - μ_0 is the prior mean.
 - σ_0^2 is the prior variance, reflecting our initial uncertainty about μ .

Goal:

- We wish to construct the posterior distribution $p(\mu|x)$, where x represents the observed data (x_1, \dots, x_n) .

Posterior for Normal Distribution Mean: Components

So we have:

1. Prior Distribution for μ : $p(\mu)$

$$p(\mu) = (2\pi\sigma_0^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 \right\}$$

Posterior for Normal Distribution Mean: Components

So we have:

1. Prior Distribution for μ : $p(\mu)$

$$p(\mu) = (2\pi\sigma_0^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2 \right\}$$

2. Likelihood for one data point x_i given μ : $p(x_i|\mu)$

$$p(x_i|\mu) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2}(x_i - \mu)^2 \right\}$$

Posterior for Normal Distribution Mean: Components

So we have:

1. Prior Distribution for μ : $p(\mu)$

$$p(\mu) = (2\pi\sigma_0^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 \right\}$$

2. Likelihood for one data point x_i given μ : $p(x_i|\mu)$

$$p(x_i|\mu) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right\}$$

3. Full Likelihood for data $x = (x_1, \dots, x_n)$ **given** μ : $p(x|\mu)$ Assuming independent observations:

$$p(x|\mu) = \prod_{i=1}^n (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right\}$$

$$p(x|\mu) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}$$

Posterior for Normal Mean: Applying Bayes' Theorem

Bayes' Theorem (proportional form):

$$p(\mu|x) \propto p(\mu) \times p(x|\mu)$$

Posterior for Normal Mean: Applying Bayes' Theorem

Bayes' Theorem (proportional form):

$$p(\mu|x) \propto p(\mu) \times p(x|\mu)$$

Substituting the Normal prior and Normal likelihood:

$$\begin{aligned} p(\mu|x) &\propto (2\pi\sigma_0^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right\} \\ &\times (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\} \end{aligned}$$

Posterior for Normal Mean: Applying Bayes' Theorem

Bayes' Theorem (proportional form):

$$p(\mu|x) \propto p(\mu) \times p(x|\mu)$$

Substituting the Normal prior and Normal likelihood:

$$\begin{aligned} p(\mu|x) &\propto (2\pi\sigma_0^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right\} \\ &\times (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\} \end{aligned}$$

We focus on terms involving μ in the exponent (dropping constants of proportionality):

$$p(\mu|x) \propto \exp\left\{-\frac{1}{2} \left[\frac{(\mu - \mu_0)^2}{\sigma_0^2} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2} \right]\right\}$$

Posterior for Normal Mean: Applying Bayes' Theorem

Bayes' Theorem (proportional form):

$$p(\mu|x) \propto p(\mu) \times p(x|\mu)$$

Substituting the Normal prior and Normal likelihood:

$$\begin{aligned} p(\mu|x) &\propto (2\pi\sigma_0^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right\} \\ &\times (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\} \end{aligned}$$

We focus on terms involving μ in the exponent (dropping constants of proportionality):

$$p(\mu|x) \propto \exp\left\{-\frac{1}{2} \left[\frac{(\mu - \mu_0)^2}{\sigma_0^2} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2} \right]\right\}$$

Posterior for Normal Mean: Applying Bayes' Theorem

Expanding the squared terms in the exponent and collecting terms in μ^2 and μ :

$$\propto \exp \left\{ -\frac{1}{2} \mu^2 \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right) + \mu \left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum x_i}{\sigma^2} \right) + \text{const} \right\}$$

This is the kernel of a Normal distribution for μ .

Posterior for Normal Distribution Mean (Continued)

General Form of a Normal PDF for a variable y with mean θ and variance ϕ :

$$f(y) = (2\pi\phi)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\phi}(y - \theta)^2 \right\}$$

This is proportional to:

$$\propto \exp \left\{ -\frac{1}{2}y^2\phi^{-1} + y\theta\phi^{-1} + \text{terms not involving } y \right\}$$

Posterior for Normal Distribution Mean (Continued)

General Form of a Normal PDF for a variable y with mean θ and variance ϕ :

$$f(y) = (2\pi\phi)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\phi}(y - \theta)^2 \right\}$$

This is proportional to:

$$\propto \exp \left\{ -\frac{1}{2}y^2\phi^{-1} + y\theta\phi^{-1} + \text{terms not involving } y \right\}$$

Comparing our posterior for μ with this general form: Our posterior was:

$$p(\mu|x) \propto \exp \left\{ -\frac{1}{2}\mu^2 \underbrace{\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)}_{\text{Posterior Precision } 1/\phi_N} + \mu \underbrace{\left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum x_i}{\sigma^2} \right)}_{\text{Term related to Post. Mean } \theta_N/\phi_N} + \text{const} \right\}$$

Posterior for Normal Distribution Mean (Continued)

General Form of a Normal PDF for a variable y with mean θ and variance ϕ :

$$f(y) = (2\pi\phi)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\phi}(y - \theta)^2 \right\}$$

This is proportional to:

$$\propto \exp \left\{ -\frac{1}{2}y^2\phi^{-1} + y\theta\phi^{-1} + \text{terms not involving } y \right\}$$

Comparing our posterior for μ with this general form: Our posterior was:

$$p(\mu|x) \propto \exp \left\{ -\frac{1}{2}\mu^2 \underbrace{\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)}_{\text{Posterior Precision } 1/\phi_N} + \mu \underbrace{\left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum x_i}{\sigma^2} \right)}_{\text{Term related to Post. Mean } \theta_N/\phi_N} + \text{const} \right\}$$

Posterior for Normal Distribution Mean (Continued)

By equating coefficients (matching terms for μ^2 and μ):

- The posterior variance ϕ_N is such that:

$$\frac{1}{\phi_N} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \quad \Rightarrow \quad \phi_N = \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}$$

- The posterior mean θ_N is such that:

$$\frac{\theta_N}{\phi_N} = \frac{\mu_0}{\sigma_0^2} + \frac{\sum x_i}{\sigma^2} \quad \Rightarrow \quad \theta_N = \phi_N \left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum x_i}{\sigma^2} \right)$$

Thus, the posterior $p(\mu|x)$ is $N(\theta_N, \phi_N)$.

Precisions and Means: Interpretation

In Bayesian statistics, **precision** (defined as $1/\text{variance}$) is often more directly additive and interpretable than variance.

For the Normal model (Normal prior for μ , Normal likelihood, known σ^2):

Precisions and Means: Interpretation

In Bayesian statistics, **precision** (defined as $1/\text{variance}$) is often more directly additive and interpretable than variance.

For the Normal model (Normal prior for μ , Normal likelihood, known σ^2):

Posterior Precision ($1/\phi_N$):

$$\frac{1}{\phi_N} = \underbrace{\frac{1}{\sigma_0^2}}_{\text{Prior Precision}} + \underbrace{\frac{n}{\sigma^2}}_{\text{Data Precision}}$$

- The posterior precision is the sum of the prior precision and the data precision.
- Data precision for the mean is n/σ^2 (increases with n , decreases with data variance σ^2).

Precisions and Means: Interpretation

In Bayesian statistics, **precision** (defined as $1/\text{variance}$) is often more directly additive and interpretable than variance.

For the Normal model (Normal prior for μ , Normal likelihood, known σ^2):

Posterior Precision ($1/\phi_N$):

$$\frac{1}{\phi_N} = \underbrace{\frac{1}{\sigma_0^2}}_{\text{Prior Precision}} + \underbrace{\frac{n}{\sigma^2}}_{\text{Data Precision}}$$

- The posterior precision is the sum of the prior precision and the data precision.
- Data precision for the mean is n/σ^2 (increases with n , decreases with data variance σ^2).

Posterior Mean (θ_N):

$$\theta_N = \frac{\left(\frac{1}{\sigma_0^2}\right) \mu_0 + \left(\frac{n}{\sigma^2}\right) \bar{x}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}}$$

- The posterior mean is a **precision-weighted average** of the prior mean (μ_0) and the data mean ($\bar{x} = \frac{\sum x_i}{n}$).
- The more precise source (prior or data) has a greater influence on the posterior mean.

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size n becomes very large?

- **Posterior Precision:**

$$\frac{1}{\phi_N} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \xrightarrow{n \rightarrow \infty} \frac{n}{\sigma^2}$$

The data precision dominates the prior precision.

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size n becomes very large?

- **Posterior Precision:**

$$\frac{1}{\phi_N} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \xrightarrow{n \rightarrow \infty} \frac{n}{\sigma^2}$$

The data precision dominates the prior precision.

- **So, Posterior Variance:**

$$\phi_N \xrightarrow{n \rightarrow \infty} \frac{\sigma^2}{n}$$

This is the familiar variance of the sample mean.

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size n becomes very large?

- **Posterior Precision:**

$$\frac{1}{\phi_N} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \xrightarrow{n \rightarrow \infty} \frac{n}{\sigma^2}$$

The data precision dominates the prior precision.

- **So, Posterior Variance:**

$$\phi_N \xrightarrow{n \rightarrow \infty} \frac{\sigma^2}{n}$$

This is the familiar variance of the sample mean.

- **Posterior Mean:**

$$\theta_N = \frac{\frac{1}{\sigma_0^2} \mu_0 + \frac{n}{\sigma^2} \bar{X}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \xrightarrow{n \rightarrow \infty} \bar{X}$$

The posterior mean converges to the sample mean (the MLE).

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size n becomes very large?

- **Posterior Precision:**

$$\frac{1}{\phi_N} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \xrightarrow{n \rightarrow \infty} \frac{n}{\sigma^2}$$

The data precision dominates the prior precision.

- **So, Posterior Variance:**

$$\phi_N \xrightarrow{n \rightarrow \infty} \frac{\sigma^2}{n}$$

This is the familiar variance of the sample mean.

- **Posterior Mean:**

$$\theta_N = \frac{\frac{1}{\sigma_0^2} \mu_0 + \frac{n}{\sigma^2} \bar{X}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \xrightarrow{n \rightarrow \infty} \bar{X}$$

The posterior mean converges to the sample mean (the MLE).

Large Sample Properties (As $n \rightarrow \infty$)

- **And so, Posterior Distribution:**

$$p(\mu|x) \xrightarrow{n \rightarrow \infty} N\left(\bar{x}, \frac{\sigma^2}{n}\right)$$

Large Sample Properties (As $n \rightarrow \infty$)

- **And so, Posterior Distribution:**

$$p(\mu|x) \xrightarrow{n \rightarrow \infty} N\left(\bar{x}, \frac{\sigma^2}{n}\right)$$

- **Comparison to Frequentist Setting:** In the frequentist setting, the sampling distribution of the sample mean \bar{x} (given true μ) is $p(\bar{x}|\mu) = N(\mu, \sigma^2/n)$.
 - The Bayesian posterior for μ and the frequentist sampling distribution for \bar{x} converge to similar forms with large n , though their interpretations differ.

Girls Heights Example: Scenario

- 10 girls aged 18 had their heights and weights measured.
- Their heights (in cm) were as follows:
 - 169.6, 166.8, 157.1, 181.1, 158.4, 165.6, 166.7, 156.5, 168.1, 165.3
- We will assume the population variance of heights is known to be $\sigma^2 = 50 \text{ cm}^2$.

Girls Heights Example: Scenario

- 10 girls aged 18 had their heights and weights measured.
- Their heights (in cm) were as follows:
 - 169.6, 166.8, 157.1, 181.1, 158.4, 165.6, 166.7, 156.5, 168.1, 165.3
- We will assume the population variance of heights is known to be $\sigma^2 = 50 \text{ cm}^2$.
- **Two individuals (Researchers/Experts) provided different prior distributions for the mean height μ :**
 - **Individual 1:** $p_1(\mu) \sim N(165, 2^2)$, i.e., $N(165, \text{variance} = 4)$
 - **Individual 2:** $p_2(\mu) \sim N(170, 3^2)$, i.e., $N(170, \text{variance} = 9)$
- We will construct the posterior for μ for each individual.

Constructing Posterior 1 (for Individual 1's Prior)

To construct the posterior, we use the formulae we have just calculated (derived for Normal-Normal model with known σ^2).

From Individual 1's Prior:

- Prior mean $\mu_0 = 165$
- Prior variance $\sigma_0^2 = 2^2 = 4$

Constructing Posterior 1 (for Individual 1's Prior)

To construct the posterior, we use the formulae we have just calculated (derived for Normal-Normal model with known σ^2).

From Individual 1's Prior:

- Prior mean $\mu_0 = 165$
- Prior variance $\sigma_0^2 = 2^2 = 4$

From the Data:

- Sample size $n = 10$
- Known population variance $\sigma^2 = 50$
- Sample mean $\bar{x} = (169.6 + \dots + 165.3)/10 = 1655.2/10 = 165.52$

Constructing Posterior 1 (for Individual 1's Prior)

To construct the posterior, we use the formulae we have just calculated (derived for Normal-Normal model with known σ^2).

From Individual 1's Prior:

- Prior mean $\mu_0 = 165$
- Prior variance $\sigma_0^2 = 2^2 = 4$

From the Data:

- Sample size $n = 10$
- Known population variance $\sigma^2 = 50$
- Sample mean $\bar{x} = (169.6 + \dots + 165.3)/10 = 1655.2/10 = 165.52$

The posterior distribution $p(\mu|x)$ will be $N(\theta_1, \phi_1)$.

Constructing Posterior 1: Girls Heights Example

- We use the formulae for posterior mean and precision (or variance) derived earlier for a Normal likelihood with a Normal prior (known data variance).
- **Prior (Individual 1):** $\mu_0 = 165$, $\sigma_0^2 = 4$ (so prior precision $1/\sigma_0^2 = 1/4 = 0.25$)
- **Data:** Heights (cm) of 10 girls.
 - Sample mean $\bar{x} = 165.52$
 - Known population variance $\sigma^2 = 50$ (so data precision per observation $1/\sigma^2 = 1/50 = 0.02$)
 - Sample size $n = 10$

Constructing Posterior 1: Calculation

The posterior distribution for μ will be $N(\theta_1, \phi_1)$.

Posterior Precision ($1/\phi_1$): Sum of prior precision and data precision for the mean (n/σ^2).

$$\frac{1}{\phi_1} = \frac{1}{4} + \frac{10}{50} = 0.25 + 0.20 = 0.45$$

So, posterior variance $\phi_1 = 1/0.45 \approx 2.222$.

Constructing Posterior 1: Calculation

The posterior distribution for μ will be $N(\theta_1, \phi_1)$.

Posterior Precision ($1/\phi_1$): Sum of prior precision and data precision for the mean (n/σ^2).

$$\frac{1}{\phi_1} = \frac{1}{4} + \frac{10}{50} = 0.25 + 0.20 = 0.45$$

So, posterior variance $\phi_1 = 1/0.45 \approx 2.222$.

Constructing Posterior 1: Girls Heights Example

Posterior Mean (θ_1): Precision-weighted average of prior mean and data mean.

$$\begin{aligned}\theta_1 &= \phi_1 \left(\frac{\mu_0}{\sigma_0^2} + \frac{n\bar{x}}{\sigma^2} \right) \\ &= \phi_1 \left(\frac{165}{4} + \frac{10 \times 165.52}{50} \right) \\ &\approx 2.222 \left(41.25 + \frac{1655.2}{50} \right) \\ &= 2.222(41.25 + 33.104) \\ &\approx 2.222 \times 74.354 \\ &\approx 165.23\end{aligned}$$

Constructing Posterior 1: Girls Heights Example

Posterior Mean (θ_1): Precision-weighted average of prior mean and data mean.

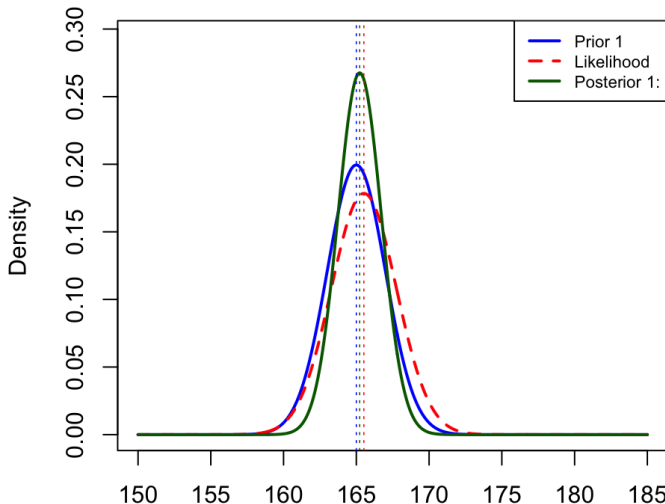
$$\begin{aligned}\theta_1 &= \phi_1 \left(\frac{\mu_0}{\sigma_0^2} + \frac{n\bar{x}}{\sigma^2} \right) \\ &= \phi_1 \left(\frac{165}{4} + \frac{10 \times 165.52}{50} \right) \\ &\approx 2.222 \left(41.25 + \frac{1655.2}{50} \right) \\ &= 2.222(41.25 + 33.104) \\ &\approx 2.222 \times 74.354 \\ &\approx 165.23\end{aligned}$$

Posterior Distribution (Individual 1):

$$p(\mu|x) \sim N(165.23, 2.222)$$

Prior and Posterior Comparison (Individual 1)

Individual 1: Prior, Likelihood, and Posterior for μ



Prior and Posterior Comparison (Individual 1)

- The prior was $N(165, 2^2)$.
- The data (likelihood) is centered around $\bar{x} = 165.52$.
- The posterior mean (165.23) is a compromise, pulled from the prior mean towards the data mean.
- The posterior variance (2.222) is smaller than both prior variance (4) and effective data variance for mean (5).

Constructing Posterior 2: Girls Heights Example

Again, to construct the posterior we use the earlier formulae.

- **Prior (Individual 2):** $\mu_0 = 170$, $\sigma_0^2 = 9$ (so prior precision $1/\sigma_0^2 = 1/9 \approx 0.111$)
- **Data (Same as before):**
 - Sample mean $\bar{x} = 165.52$
 - Known population variance $\sigma^2 = 50$
 - Sample size $n = 10$

The posterior distribution for μ will be $N(\theta_2, \phi_2)$.

Constructing Posterior 2: Calculation

Posterior Precision ($1/\phi_2$):

$$\frac{1}{\phi_2} = \frac{1}{9} + \frac{10}{50} \approx 0.1111 + 0.20 = 0.3111$$

So, posterior variance $\phi_2 = 1/0.3111 \approx 3.214$.

Constructing Posterior 2: Calculation

Posterior Precision ($1/\phi_2$):

$$\frac{1}{\phi_2} = \frac{1}{9} + \frac{10}{50} \approx 0.1111 + 0.20 = 0.3111$$

So, posterior variance $\phi_2 = 1/0.3111 \approx 3.214$.

Posterior Mean (θ_2):

$$\theta_2 = \phi_2 \left(\frac{\mu_0}{\sigma_0^2} + \frac{n\bar{x}}{\sigma^2} \right) = \phi_2 \left(\frac{170}{9} + \frac{10 \times 165.52}{50} \right)$$

$$\theta_2 \approx 3.214 (18.8889 + 33.104) = 3.214 \times 51.9929 \approx 167.12$$

Constructing Posterior 2: Calculation

Posterior Precision ($1/\phi_2$):

$$\frac{1}{\phi_2} = \frac{1}{9} + \frac{10}{50} \approx 0.1111 + 0.20 = 0.3111$$

So, posterior variance $\phi_2 = 1/0.3111 \approx 3.214$.

Posterior Mean (θ_2):

$$\theta_2 = \phi_2 \left(\frac{\mu_0}{\sigma_0^2} + \frac{n\bar{x}}{\sigma^2} \right) = \phi_2 \left(\frac{170}{9} + \frac{10 \times 165.52}{50} \right)$$

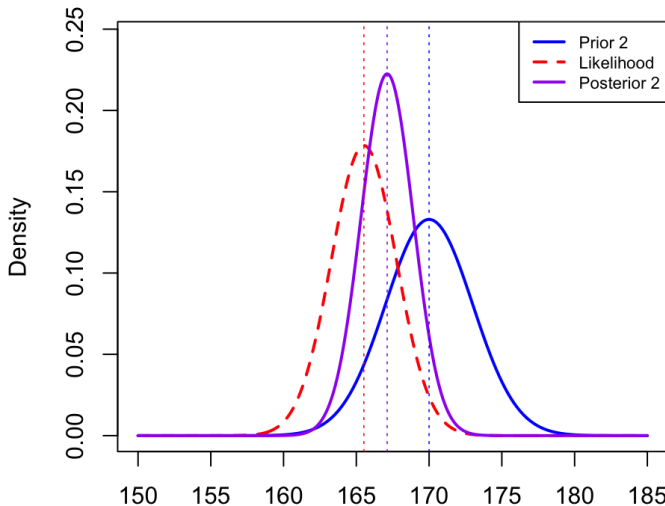
$$\theta_2 \approx 3.214 (18.8889 + 33.104) = 3.214 \times 51.9929 \approx 167.12$$

Posterior Distribution (Individual 2):

$$p(\mu|x) \sim N(167.12, 3.214)$$

Prior 2 Comparison

Individual 2: Prior, Likelihood, and Posterior for μ



Prior 2 Comparison

- Note: This prior ($N(170, 3^2)$) is not as close to the data mean ($\bar{x} = 165.52$) as prior 1 was.
- The posterior mean (167.12) is still a compromise, pulled from the prior mean (170) towards the data mean (165.52).
- The posterior is somewhere between the prior and the likelihood, reflecting the influence of both.

Other Conjugate Examples

When the posterior distribution is in the same probability distribution family as the prior distribution, we have **conjugacy**. This often simplifies calculations.

Other Conjugate Examples

When the posterior distribution is in the same probability distribution family as the prior distribution, we have **conjugacy**. This often simplifies calculations. **Examples include:**

Likelihood	Parameter	Prior	Posterior
Normal	Mean (known variance)	Normal	Normal
Normal	Precision (known mean)	Gamma	Gamma
Binomial	Probability p	Beta	Beta
Poisson	Mean λ	Gamma	Gamma

Table: Common Conjugate Prior-Likelihood Pairs.

- Conjugacy was common before widespread MCMC methods, as it allowed for analytical solutions.
- With MCMC, conjugacy is less critical but can still be convenient.

In All (Conjugate) Cases: General Properties

- The **posterior mean** is a compromise (often a weighted average) between the prior mean and the information from the data (e.g., sample mean or Maximum Likelihood Estimate - MLE).

In All (Conjugate) Cases: General Properties

- The **posterior mean** is a compromise (often a weighted average) between the prior mean and the information from the data (e.g., sample mean or Maximum Likelihood Estimate - MLE).
- The **posterior standard deviation** (or variance/precision) is typically less than both the prior s.d. and the standard error from the data alone (s.e. of MLE).
 - This reflects that combining information sources increases precision (reduces uncertainty).

In All (Conjugate) Cases: General Properties

- The **posterior mean** is a compromise (often a weighted average) between the prior mean and the information from the data (e.g., sample mean or Maximum Likelihood Estimate - MLE).
- The **posterior standard deviation** (or variance/precision) is typically less than both the prior s.d. and the standard error from the data alone (s.e. of MLE).
 - This reflects that combining information sources increases precision (reduces uncertainty).
- **Quote (Senn):**
'A Bayesian is one who, vaguely expecting a horse and catching a glimpse of a donkey, strongly concludes he has seen a mule'
(Illustrates the "compromise" nature of the posterior).

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size (n) becomes very large?

- The **posterior mean** \rightarrow the MLE (Maximum Likelihood Estimate).
 - The data "overwhelms" the prior.

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size (n) becomes very large?

- The **posterior mean** \rightarrow the MLE (Maximum Likelihood Estimate).
 - The data "overwhelms" the prior.
- The **posterior standard deviation** \rightarrow the s.e. (MLE).
 - Uncertainty becomes dominated by the sampling variability in the data.

Large Sample Properties (As $n \rightarrow \infty$)

What happens as the sample size (n) becomes very large?

- The **posterior mean** \rightarrow the MLE (Maximum Likelihood Estimate).
 - The data "overwhelms" the prior.
- The **posterior standard deviation** \rightarrow the s.e. (MLE).
 - Uncertainty becomes dominated by the sampling variability in the data.
- The **posterior distribution does not depend (much) on the prior**.
 - (Assuming the prior was not dogmatic, i.e., didn't assign zero probability to plausible parameter values).

Non-informative Priors

- We often do not have strong prior information, or wish to let the data "speak for itself" as much as possible.
 - (Though "true Bayesians" might argue we *always* have some implicit prior information!)
- In such cases, we aim for good agreement between the frequentist approach and the Bayesian approach using a so-called "non-informative" prior.

Non-informative Priors

- We often do not have strong prior information, or wish to let the data "speak for itself" as much as possible.
 - (Though "true Bayesians" might argue we *always* have some implicit prior information!)
- In such cases, we aim for good agreement between the frequentist approach and the Bayesian approach using a so-called "non-informative" prior.
- **Better terms:** "Diffuse" or "flat" priors, as no prior is strictly non-informative (e.g., uniform on one scale is not uniform on a transformed scale).

Non-informative Priors

- We often do not have strong prior information, or wish to let the data "speak for itself" as much as possible.
 - (Though "true Bayesians" might argue we *always* have some implicit prior information!)
- In such cases, we aim for good agreement between the frequentist approach and the Bayesian approach using a so-called "non-informative" prior.
- **Better terms:** "Diffuse" or "flat" priors, as no prior is strictly non-informative (e.g., uniform on one scale is not uniform on a transformed scale).
- **For our example of an unknown mean (Normal likelihood):**
 - A Uniform distribution over a very large range.
 - A Normal distribution with a huge variance (e.g., $N(0, 10000^2)$).
- Goal: To have minimal influence from the prior on the posterior, especially with reasonable amounts of data.

Point and Interval Estimation from the Posterior

- In Bayesian inference, the primary outcome of interest for a parameter is its **full posterior distribution**.
- However, we are often interested in summaries of this distribution.

Point and Interval Estimation from the Posterior

- In Bayesian inference, the primary outcome of interest for a parameter is its **full posterior distribution**.
- However, we are often interested in summaries of this distribution.
- **Point Estimate:**
 - A simple point estimate could be the **mean** of the posterior.
 - Alternatives include the **median** or **mode** of the posterior.

Point and Interval Estimation from the Posterior

- In Bayesian inference, the primary outcome of interest for a parameter is its **full posterior distribution**.
- However, we are often interested in summaries of this distribution.
- **Point Estimate:**
 - A simple point estimate could be the **mean** of the posterior.
 - Alternatives include the **median** or **mode** of the posterior.
- **Interval Estimates:**
 - Easy to obtain from the posterior distribution.
 - Given several names, all referring to essentially the same concept:
 - **Credible Intervals** (most common Bayesian term)
 - Bayesian confidence intervals
 - Highest Density Regions (HDR) or Highest Posterior Density (HPD) intervals

Credible Intervals: Heights Example

Recall Posterior for Individual 1 (Prior $N(165, 2^2)$):

$$P(\mu|x) \sim N(165.23, \text{variance} = 2.222)$$

So, posterior standard deviation $\sqrt{2.222} \approx 1.49$.

Credible Intervals: Heights Example

Recall Posterior for Individual 1 (Prior $N(165, 2^2)$):

$$P(\mu|x) \sim N(165.23, \text{variance} = 2.222)$$

So, posterior standard deviation $\sqrt{2.222} \approx 1.49$. A 95% credible interval for μ is (using mean $\pm 1.96 \times$ SD for Normal):

$$165.23 \pm 1.96 \times 1.4907 \approx 165.23 \pm 2.92$$

Which gives (162.31, 168.15).

Credible Intervals: Heights Example

Recall Posterior for Individual 1 (Prior $N(165, 2^2)$):

$$P(\mu|x) \sim N(165.23, \text{variance} = 2.222)$$

So, posterior standard deviation $\sqrt{2.222} \approx 1.49$. A 95% credible interval for μ is (using mean $\pm 1.96 \times$ SD for Normal):

$$165.23 \pm 1.96 \times 1.4907 \approx 165.23 \pm 2.92$$

Which gives (162.31, 168.15).

Recall Posterior for Individual 2 (Prior $N(170, 3^2)$):

$$P(\mu|x) \sim N(167.12, \text{variance} = 3.214)$$

A 95% credible interval for μ is (using mean $\pm 1.96 \times$ SD): Which gives (163.61, 170.63).

Credible Interval Interpretation

Key Distinction

Credible intervals can be interpreted in the more natural way:

- For a 95% credible interval, e.g., (162.31, 168.15) for μ :
 - "There is a 95% probability that the true mean height μ lies within the interval (162.31, 168.15)."

Credible Interval Interpretation

Key Distinction

Credible intervals can be interpreted in the more natural way:

- For a 95% credible interval, e.g., (162.31, 168.15) for μ :
 - "There is a 95% probability that the true mean height μ lies within the interval (162.31, 168.15)."
- This is different from the frequentist interpretation of a confidence interval, which states that:
 - "95% of such similarly constructed confidence intervals (from repeated sampling) would contain the true μ ." (The probability refers to the procedure, not the specific interval).

The Bayesian interpretation is often what people intuitively (but often incorrectly) ascribe to frequentist confidence intervals.

Outline for this Section

- 1 Introduction and Session Aims (9:30 – 10:00)
- 2 Theoretical – Bayes (Part 1) (10:00 – 11:30)
 - What is Bayesian Statistics?
 - Statistical Distributions
 - Priors
 - Example: Normal Data, Unknown Mean (Known Variance)
- 3 Theoretical – Bayes (Part 2) (11:50am – 12:10pm)
 - The Computational Challenge Intro to MCMC
- 4 Summary of Morning Session (12:10 – 12:30pm)
- 5 Practical Application (1:30 – 4:15pm)
 - Tools (Stan Overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison
- 6 Review and Closing (4:15 – 4:30pm)

The Computational Challenge

Posterior: $P(\theta|y) \propto P(y|\theta)P(\theta)$.

The denominator $P(y) = \int P(y|\theta')P(\theta')d\theta'$ is often intractable.

The Computational Challenge

Posterior: $P(\theta|y) \propto P(y|\theta)P(\theta)$.

The denominator $P(y) = \int P(y|\theta')P(\theta')d\theta'$ is often intractable.

Solution: Sample from the posterior using MCMC.

Algorithms to draw samples from a target posterior distribution.

- Construct a Markov chain whose stationary distribution is $P(\theta|y)$.
- After "burn-in", samples approximate the posterior.

Goal: Sample from joint posterior $p(\beta_0, \beta_1, \sigma^2|y)$ for a linear model.

Gibbs Sampling

An MCMC algorithm used when full conditional posteriors are known.
Iteratively sample each parameter conditional on current values of others:

- ① Initialize $\theta_1^{(0)}, \dots, \theta_k^{(0)}$.
- ② For $t = 1, \dots, T$:
 - Sample $\theta_1^{(t)} \sim p(\theta_1 | \theta_2^{(t-1)}, \dots, \theta_k^{(t-1)}, y)$
 - Sample $\theta_2^{(t)} \sim p(\theta_2 | \theta_1^{(t)}, \theta_3^{(t-1)}, \dots, y)$
 - ...
 - Sample $\theta_k^{(t)} \sim p(\theta_k | \theta_1^{(t)}, \dots, \theta_{k-1}^{(t)}, y)$

Computationally efficient if conditionals are standard (e.g. Normal, Gamma).

Other MCMC Samplers (Metropolis-Hastings, HMC/NUTS)

- **Metropolis-Hastings:** More general. Propose accept/reject. Used if analytical (conditional) posterior unknown.
- **HMC/NUTS (Stan):** Uses gradient information for efficient proposals. Excellent for many continuous parameter problems.

Why HMC? Limitations of Simpler MCMC:

- Basic MCMC (e.g., Random Walk Metropolis) explores the parameter space by taking small, random steps.
- This can be very inefficient for:
 - High-dimensional posteriors.
 - Correlated parameters.
- Leads to slow mixing, high autocorrelation, and many rejected proposals.

Why HMC? Limitations of Simpler MCMC:

- Basic MCMC (e.g., Random Walk Metropolis) explores the parameter space by taking small, random steps.
- This can be very inefficient for:
 - High-dimensional posteriors.
 - Correlated parameters.
- Leads to slow mixing, high autocorrelation, and many rejected proposals.

HMC: Smarter Proposals using Physics Analogy

- HMC introduces auxiliary "momentum" variables for each parameter.
- It simulates the dynamics of a particle moving over a surface defined by the (negative log) posterior density (the "potential energy").
- This allows for **longer, more directed moves** to distant, high-probability regions of the posterior.
- A Metropolis-Hastings acceptance step corrects for numerical errors in simulating the trajectory, ensuring convergence to the true posterior, typically with very high acceptance rates.

Key Advantages

More efficient exploration, lower autocorrelation, better for complex models. (Core of Stan's sampler)

Challenge in Standard HMC:

- Two key tuning parameters need to be set by the user:
 - ① **Step size (ϵ):** How large are the steps in the trajectory simulation?
 - ② **Number of steps (L):** How long is the trajectory simulated for?
- Poor choices for ϵ and L can lead to inefficiency or incorrect sampling. Optimal tuning can be difficult.

Challenge in Standard HMC:

- Two key tuning parameters need to be set by the user:
 - ① **Step size (ϵ):** How large are the steps in the trajectory simulation?
 - ② **Number of steps (L):** How long is the trajectory simulated for?
- Poor choices for ϵ and L can lead to inefficiency or incorrect sampling. Optimal tuning can be difficult.

No-U-Turn Sampler (NUTS): Automating HMC

NUTS: Adaptive HMC (Default in Stan)

- NUTS is an extension of HMC that **automates the selection of L** .
- It builds a path (trajectory) by taking steps, and adaptively stops when the path starts to "U-turn" (i.e., move back towards where it started).
 - This prevents overshooting good regions or taking too short trajectories.
- NUTS also incorporates sophisticated algorithms to automatically adapt the step size ϵ during the warmup/adaptation phase of MCMC.

No-U-Turn Sampler (NUTS): Automating HMC

NUTS: Adaptive HMC (Default in Stan)

- NUTS is an extension of HMC that **automates the selection of L** .
- It builds a path (trajectory) by taking steps, and adaptively stops when the path starts to "U-turn" (i.e., move back towards where it started).
 - This prevents overshooting good regions or taking too short trajectories.
- NUTS also incorporates sophisticated algorithms to automatically adapt the step size ϵ during the warmup/adaptation phase of MCMC.

Why NUTS is Powerful

Significantly reduces the need for manual sampler tuning, making HMC robust and efficient for a wide range of models. This is a major reason for Stan's effectiveness.

How Does One Fit Models in a Bayesian Framework?

How Does One Fit Models in a Bayesian Framework?

- In the first section, we illustrated using **conjugate priors** to analytically evaluate a posterior distribution for a model with one unknown parameter.
- This is often not feasible for more complex models.

How Does One Fit Models in a Bayesian Framework?

- In the first section, we illustrated using **conjugate priors** to analytically evaluate a posterior distribution for a model with one unknown parameter.
- This is often not feasible for more complex models.
- Let us now consider a simple linear regression:

$$\text{weight}_i = \beta_0 + \beta_1 \text{height}_i + e_i$$

$$e_i \sim N(0, \sigma^2)$$

How Does One Fit Models in a Bayesian Framework?

- In the first section, we illustrated using **conjugate priors** to analytically evaluate a posterior distribution for a model with one unknown parameter.
- This is often not feasible for more complex models.
- Let us now consider a simple linear regression:

$$\text{weight}_i = \beta_0 + \beta_1 \text{height}_i + e_i$$

$$e_i \sim N(0, \sigma^2)$$

- If we used conjugate priors (for illustration of a past method):

$$\beta_0 \sim N(0, m_0), \quad \beta_1 \sim N(0, m_1)$$

$$\sigma^2 \sim \text{Inverse-Gamma}(\varepsilon, \varepsilon)$$

(Example: $m_0 = m_1 = 10^6, \varepsilon = 10^{-3}$ for non-informative priors)

How Does One Fit Models in a Bayesian Framework?

- In the first section, we illustrated using **conjugate priors** to analytically evaluate a posterior distribution for a model with one unknown parameter.
- This is often not feasible for more complex models.
- Let us now consider a simple linear regression:

$$\text{weight}_i = \beta_0 + \beta_1 \text{height}_i + e_i$$

$$e_i \sim N(0, \sigma^2)$$

- If we used conjugate priors (for illustration of a past method):

$$\beta_0 \sim N(0, m_0), \quad \beta_1 \sim N(0, m_1)$$

$$\sigma^2 \sim \text{Inverse-Gamma}(\varepsilon, \varepsilon)$$

(Example: $m_0 = m_1 = 10^6, \varepsilon = 10^{-3}$ for non-informative priors)

Our Goal (Generally): Make inferences on the joint posterior distribution:

$$p(\beta_0, \beta_1, \sigma^2 | \text{data})$$

For complex models, this joint posterior is hard to obtain directly.
This is where MCMC methods come in.

Goal:

- To sample from the joint posterior distribution:

$$p(\beta_0, \beta_1, \sigma^2 | \text{data})$$

Goal:

- To sample from the joint posterior distribution:

$$p(\beta_0, \beta_1, \sigma^2 | \text{data})$$

Problem:

- For complex models, directly calculating or sampling from this joint posterior often involves high-dimensional integration, which is analytically or computationally intractable.

Goal:

- To sample from the joint posterior distribution:

$$p(\beta_0, \beta_1, \sigma^2 | \text{data})$$

Problem:

- For complex models, directly calculating or sampling from this joint posterior often involves high-dimensional integration, which is analytically or computationally intractable.

Goal: Solution (MCMC Approach):

- It may be possible to sample from simpler **conditional posterior distributions**.
- For example, for parameters $\beta_0, \beta_1, \sigma^2$:
 - $p(\beta_0|\text{data}, \beta_1, \sigma^2)$ (sample β_0 given data and current β_1, σ^2)
 - $p(\beta_1|\text{data}, \beta_0, \sigma^2)$ (sample β_1 given data and current β_0, σ^2)
 - $p(\sigma^2|\text{data}, \beta_0, \beta_1)$ (sample σ^2 given data and current β_0, β_1)
- MCMC algorithms iteratively sample from these conditionals.
- It can be shown that after *convergence*, such a sampling approach generates (dependent) samples from the target *joint* posterior distribution.

Gibbs Sampling: A Type of MCMC

- When we can sample *directly* from the full conditional posterior distributions for each parameter (or block of parameters), the algorithm is known as **Gibbs Sampling**.

Gibbs Sampling: A Type of MCMC

- When we can sample *directly* from the full conditional posterior distributions for each parameter (or block of parameters), the algorithm is known as **Gibbs Sampling**.
- **Procedure for the linear regression example $(\beta_0, \beta_1, \sigma^2)$:**

- 1 **Initialization:** Give all unknown parameters starting values:

$$\beta_0^{(0)}, \beta_1^{(0)}, (\sigma^2)^{(0)}$$

- 2 **Iteration Loop (for $t = 1, 2, \dots, T$):** Sequentially sample each parameter from its full conditional distribution, using the most recently updated values of the other parameters.

Gibbs Sampling: Iterative Steps (Example)

At iteration t :

1 Sample $\beta_0^{(t)}$ from $p(\beta_0|\text{data}, \beta_1^{(t-1)}, (\sigma^2)^{(t-1)})$

- This generates $\beta_0^{(t)}$.

Gibbs Sampling: Iterative Steps (Example)

At iteration t :

- 1 Sample $\beta_0^{(t)}$ from $p(\beta_0|\text{data}, \beta_1^{(t-1)}, (\sigma^2)^{(t-1)})$
 - This generates $\beta_0^{(t)}$.
- 2 Sample $\beta_1^{(t)}$ from $p(\beta_1|\text{data}, \beta_0^{(t)}, (\sigma^2)^{(t-1)})$
 - Note: Uses the just-updated $\beta_0^{(t)}$.
 - This generates $\beta_1^{(t)}$.

Gibbs Sampling: Iterative Steps (Example)

At iteration t :

- 1 Sample $\beta_0^{(t)}$ from $p(\beta_0|\text{data}, \beta_1^{(t-1)}, (\sigma^2)^{(t-1)})$
 - This generates $\beta_0^{(t)}$.
- 2 Sample $\beta_1^{(t)}$ from $p(\beta_1|\text{data}, \beta_0^{(t)}, (\sigma^2)^{(t-1)})$
 - Note: Uses the just-updated $\beta_0^{(t)}$.
 - This generates $\beta_1^{(t)}$.
- 3 Sample $(\sigma^2)^{(t)}$ from $p(\sigma^2|\text{data}, \beta_0^{(t)}, \beta_1^{(t)})$
 - Uses updated $\beta_0^{(t)}$ and $\beta_1^{(t)}$.
 - This generates $(\sigma^2)^{(t)}$.

Gibbs Sampling: Iterative Steps (Example)

At iteration t :

- 1 Sample $\beta_0^{(t)}$ from $p(\beta_0|\text{data}, \beta_1^{(t-1)}, (\sigma^2)^{(t-1)})$
 - This generates $\beta_0^{(t)}$.
- 2 Sample $\beta_1^{(t)}$ from $p(\beta_1|\text{data}, \beta_0^{(t)}, (\sigma^2)^{(t-1)})$
 - Note: Uses the just-updated $\beta_0^{(t)}$.
 - This generates $\beta_1^{(t)}$.
- 3 Sample $(\sigma^2)^{(t)}$ from $p(\sigma^2|\text{data}, \beta_0^{(t)}, \beta_1^{(t)})$
 - Uses updated $\beta_0^{(t)}$ and $\beta_1^{(t)}$.
 - This generates $(\sigma^2)^{(t)}$.

Gibbs Sampling: Iterative Steps (Example)

Process:

- These steps are repeated for many iterations.
- The sequence of sampled values $\{(\beta_0^{(t)}, \beta_1^{(t)}, (\sigma^2)^{(t)})\}_{t=1}^T$ forms a **Markov chain**.
- It is hoped (and under certain conditions, guaranteed) that this chain converges to its equilibrium distribution, which is the target joint posterior distribution $p(\beta_0, \beta_1, \sigma^2 | \text{data})$.

Calculating the Conditional Distributions (for Gibbs)

- For the Gibbs sampling algorithm to work, we need to be able to sample from the conditional posterior distributions.
- If these conditional distributions have standard forms (e.g., Normal, Gamma, Beta), then it is (relatively) easy to draw random samples from them.

Calculating the Conditional Distributions (for Gibbs)

- For the Gibbs sampling algorithm to work, we need to be able to sample from the conditional posterior distributions.
- If these conditional distributions have standard forms (e.g., Normal, Gamma, Beta), then it is (relatively) easy to draw random samples from them.
- **How are they derived? (Conceptual):**
 - 1 Write down the full joint posterior (proportional to prior \times likelihood).
 - 2 To find $p(\theta_j | \text{data}, \theta_{-j})$ (conditional for θ_j given all other parameters θ_{-j}):
 - Treat all parameters θ_{-j} and data as constants.
 - Identify terms involving θ_j .
 - Try to match the resulting functional form (kernel) to a known standard distribution.

Calculating the Conditional Distributions (for Gibbs)

- For the Gibbs sampling algorithm to work, we need to be able to sample from the conditional posterior distributions.
- If these conditional distributions have standard forms (e.g., Normal, Gamma, Beta), then it is (relatively) easy to draw random samples from them.
- **How are they derived? (Conceptual):**
 - 1 Write down the full joint posterior (proportional to prior \times likelihood).
 - 2 To find $p(\theta_j | \text{data}, \theta_{-j})$ (conditional for θ_j given all other parameters θ_{-j}):
 - Treat all parameters θ_{-j} and data as constants.
 - Identify terms involving θ_j .
 - Try to match the resulting functional form (kernel) to a known standard distribution.

Matching Distributional Forms (Reference)

Key Idea: Recognize the "kernel" of a distribution.

Normal Distribution: If a parameter θ follows a Normal(μ, σ^2) distribution, then its PDF $p(\theta)$ is proportional to:

$$p(\theta) \propto \exp(a\theta^2 + b\theta)$$

where $a = -\frac{1}{2\sigma^2}$ and $b = \frac{\mu}{\sigma^2}$.

Gamma Distribution: If a parameter θ follows a Gamma(α, β) distribution (shape α , rate β), its PDF $p(\theta)$ is proportional to:

$$p(\theta) \propto \theta^{\alpha-1} \exp(-\beta\theta)$$

Here, if we see $p(\theta) \propto \theta^A \exp(B\theta)$, then $A = \alpha - 1$ and $B = -\beta$. (Note: *Parameterizations of Gamma can vary. Be consistent!*)

Gibbs Sampling Algorithm Summary (Linear Regression)

Repeat the following three steps iteratively:

- 1 Generate β_0 from its Normal conditional posterior distribution (given current β_1, σ^2).
- 2 Generate β_1 from its Normal conditional posterior distribution (given current β_0, σ^2).
- 3 Generate $1/\sigma^2$ (precision) from its Gamma conditional posterior distribution (given current β_0, β_1). (Then $\sigma^2 = 1/\text{precision}$).

Gibbs Sampling Algorithm Summary (Linear Regression)

Repeat the following three steps iteratively:

- 1 Generate β_0 from its Normal conditional posterior distribution (given current β_1, σ^2).
- 2 Generate β_1 from its Normal conditional posterior distribution (given current β_0, σ^2).
- 3 Generate $1/\sigma^2$ (precision) from its Gamma conditional posterior distribution (given current β_0, β_1). (Then $\sigma^2 = 1/\text{precision}$).

Gibbs Sampling Algorithm Summary (Linear Regression)

Convergence and Burn-in: Two critical questions immediately arise:

- 1 **Starting Values** **Convergence:** We start from arbitrary starting values. When can we safely say that our samples are actually from the correct target (posterior) distribution? (This relates to the "burn-in" period).

Gibbs Sampling Algorithm Summary (Linear Regression)

Convergence and Burn-in: Two critical questions immediately arise:

- 1 **Starting Values Convergence:** We start from arbitrary starting values. When can we safely say that our samples are actually from the correct target (posterior) distribution? (This relates to the "burn-in" period).
- 2 **Chain Length:** After this burn-in point (convergence), how long should we run the chain for and store values to get reliable estimates of the posterior?

Checking Convergence: The Researcher's Responsibility!

Checking Convergence: The Researcher's Responsibility!

- **Responsibility:** Assessing convergence is crucial and falls on the researcher.
- **Target:** Convergence is to a target *distribution* (the required posterior), not to a single point estimate like in Maximum Likelihood (ML) methods.

Checking Convergence: The Researcher's Responsibility!

- **Responsibility:** Assessing convergence is crucial and falls on the researcher.
- **Target:** Convergence is to a target *distribution* (the required posterior), not to a single point estimate like in Maximum Likelihood (ML) methods.
- **Visual Indication:** Once convergence has been reached (after burn-in), samples from the chain should look like a random scatter around a stable mean value. This is often assessed visually using **trace plots**.

Checking Convergence: The Researcher's Responsibility!

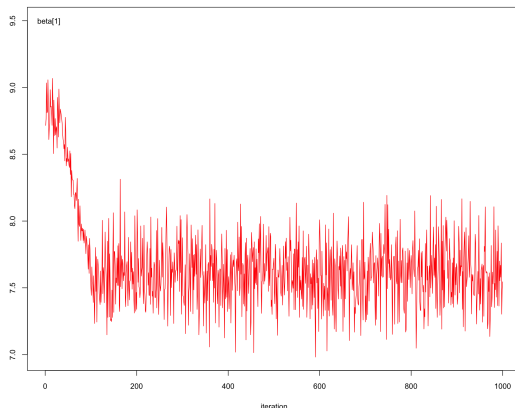


Figure: Example Trace Plot for a parameter (e.g., $\beta[1]$). Convergence appears to occur around 100 iterations in this example, followed by stable sampling.

How Many Iterations After Convergence?

- After the initial "burn-in" period (where chains converge), further iterations are needed to obtain a sufficient number of samples for posterior inference.
- **More iterations = more accurate posterior estimates** (e.g., mean, quantiles).

How Many Iterations After Convergence?

- After the initial "burn-in" period (where chains converge), further iterations are needed to obtain a sufficient number of samples for posterior inference.
- **More iterations = more accurate posterior estimates** (e.g., mean, quantiles).
- **Dependence in Chains:** MCMC chains produce *dependent* samples (autocorrelated). The degree of dependence or autocorrelation in the chain will influence how many iterations are needed for a given level of accuracy.
 - High autocorrelation means samples are very similar to previous ones, so more samples are needed to get the same amount of "information" as independent samples.

How Many Iterations After Convergence?

- After the initial "burn-in" period (where chains converge), further iterations are needed to obtain a sufficient number of samples for posterior inference.
- **More iterations = more accurate posterior estimates** (e.g., mean, quantiles).
- **Dependence in Chains:** MCMC chains produce *dependent* samples (autocorrelated). The degree of dependence or autocorrelation in the chain will influence how many iterations are needed for a given level of accuracy.
 - High autocorrelation means samples are very similar to previous ones, so more samples are needed to get the same amount of "information" as independent samples.

How Many Iterations After Convergence?

- **Assessing Accuracy:** The accuracy of posterior estimates (e.g., posterior mean) can be assessed by the **Monte Carlo Standard Error (MCSE)** for each parameter.
 - MCSE quantifies the uncertainty in the estimate due to the MCMC sampling process (not to be confused with posterior standard deviation).

Diagnostic: Trace Plot

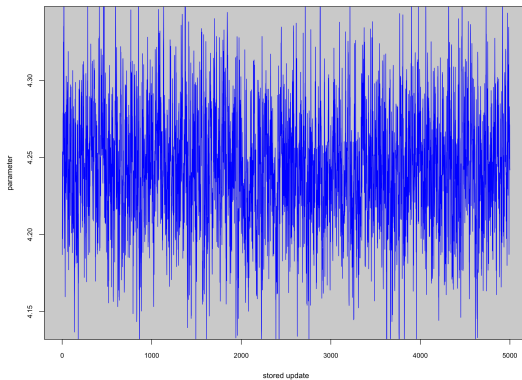


Figure: Example Trace Plot.

- This graph plots the generated values of a parameter against the iteration number.
- **What to look for:**
 - After burn-in, the chain should look like a "fat hairy caterpillar" – stationary, with good mixing around a stable mean.
 - No long-term trends or drifts.
 - No persistent periodicities.
- A crude (and informal) test of mixing is the "blue finger" test (if the plot is dense and looks like a solid band).
- The example chain "doesn't mix that well but could be worse!" (Suggests some autocorrelation might be present).
- Running multiple chains from different starting points and overlaying their trace plots is highly recommended.

Diagnostic: Kernel Density Plot

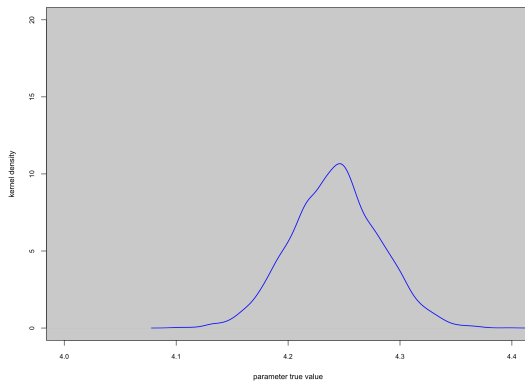


Figure: Example Kernel Density Plot.

- This plot is like a smoothed histogram of the posterior samples (after burn-in).
- Instead of counting estimates into discrete bins (like a histogram), the effect of each iteration's sample value is spread around that value via a Kernel function.
- The density at each point is the sum of these kernel function contributions from all post-burn-in samples.
- Provides a visual representation of the shape of the marginal posterior distribution for a parameter.
- The Kernel density plot has a "smoothness" parameter (bandwidth) that can be modified, affecting its appearance.

• Autocorrelation Function (ACF):

- Measures how correlated the values in the chain are with their "neighbors" at different lags (distances).
- Lag k is the correlation between X_t and X_{t-k} .
- An independent chain (ideal but rare with MCMC) would have approximately zero autocorrelation at all lags (except lag 0).
- High ACF values that decay slowly indicate poor mixing and high dependence between samples.

• Partial Autocorrelation Function (PACF):

- Measures the correlation between X_t and X_{t-k} after removing the effect of the intermediate lags $(X_{t-1}, \dots, X_{t-k+1})$.
- For an AR(1) process (where X_t only depends on X_{t-1}), the PACF should cut off after lag 1.
- In MCMC, we want ACF to drop to near zero relatively quickly.

Diagnostic: Time Series Plots (ACF and PACF)

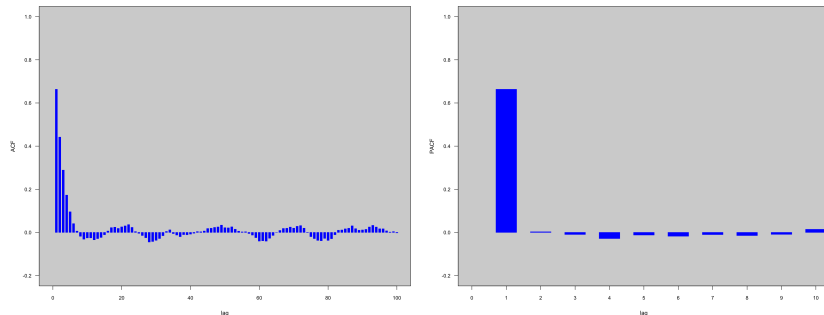


Figure: Example Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) plots.

Once we have a (converged) chain of posterior samples, we calculate summary statistics:

- **Three common estimates of location (central tendency):**

- **Mean:** Arithmetic average of the post-burn-in samples from the chain.
- **Mode:** Peak of the (kernel) density plot; most probable value.
- **Median (50% quantile):** Middle value after sorting the chain samples. Robust to outliers.

Once we have a (converged) chain of posterior samples, we calculate summary statistics:

- **Three common estimates of location (central tendency):**

- **Mean:** Arithmetic average of the post-burn-in samples from the chain.
- **Mode:** Peak of the (kernel) density plot; most probable value.
- **Median (50% quantile):** Middle value after sorting the chain samples.
Robust to outliers.

- **Spread/Uncertainty:**

- **Standard Deviation (SD):** Calculated from the chain samples, represents posterior uncertainty.
- **Quantiles (e.g., 2.5%, 97.5%):** Used to form credible intervals (possibly non-symmetric).

Once we have a (converged) chain of posterior samples, we calculate summary statistics:

- **Three common estimates of location (central tendency):**

- **Mean:** Arithmetic average of the post-burn-in samples from the chain.
- **Mode:** Peak of the (kernel) density plot; most probable value.
- **Median (50% quantile):** Middle value after sorting the chain samples.
Robust to outliers.

- **Spread/Uncertainty:**

- **Standard Deviation (SD):** Calculated from the chain samples, represents posterior uncertainty.
- **Quantiles (e.g., 2.5%, 97.5%):** Used to form credible intervals (possibly non-symmetric).

- The Monte Carlo Standard Error is an indication of how much numerical error is in the posterior estimate (e.g., posterior mean) due to the fact that MCMC is a stochastic simulation method.
- It tells us about the precision of our MCMC-based estimate of a posterior quantity.

- The Monte Carlo Standard Error is an indication of how much numerical error is in the posterior estimate (e.g., posterior mean) due to the fact that MCMC is a stochastic simulation method.
- It tells us about the precision of our MCMC-based estimate of a posterior quantity.
- As the number of (post-burn-in) iterations increases, the MCSE $\rightarrow 0$.

- The Monte Carlo Standard Error is an indication of how much numerical error is in the posterior estimate (e.g., posterior mean) due to the fact that MCMC is a stochastic simulation method.
- It tells us about the precision of our MCMC-based estimate of a posterior quantity.
- As the number of (post-burn-in) iterations increases, the MCSE $\rightarrow 0$.
- For an *independent* sampler (not typical MCMC), MCSE for the mean would be SD/\sqrt{n} (where SD is posterior SD, n is number of samples).
- However, for MCMC, MCSE is adjusted due to the autocorrelation in the chain (often larger than SD/\sqrt{n} for the same n).
- **Rule of thumb:** MCSE should be small relative to the posterior SD (e.g., $MCSE \leq 5\% \text{ of posterior SD}$).

Effective Sample Size (ESS)

- This quantity gives an estimate of the equivalent number of *independent* iterations that the MCMC chain represents.
- MCMC samples are autocorrelated, so N MCMC samples typically contain less information than N truly independent samples from the posterior.

Effective Sample Size (ESS)

- This quantity gives an estimate of the equivalent number of *independent* iterations that the MCMC chain represents.
- MCMC samples are autocorrelated, so N MCMC samples typically contain less information than N truly independent samples from the posterior.
- ESS is related to the Autocorrelation Function (ACF) and the MCSE.
- Its formula is (conceptually):

$$\text{ESS} = n/\kappa \quad \text{where } \kappa = 1 + 2 \sum_{k=1}^{\infty} \rho(k)$$

(n is actual number of post-burn-in samples, $\rho(k)$ is autocorrelation at lag k , κ is integrated autocorrelation time).

Effective Sample Size (ESS)

- This quantity gives an estimate of the equivalent number of *independent* iterations that the MCMC chain represents.
- MCMC samples are autocorrelated, so N MCMC samples typically contain less information than N truly independent samples from the posterior.
- ESS is related to the Autocorrelation Function (ACF) and the MCSE.
- Its formula is (conceptually):

$$\text{ESS} = n/\kappa \quad \text{where } \kappa = 1 + 2 \sum_{k=1}^{\infty} \rho(k)$$

(n is actual number of post-burn-in samples, $\rho(k)$ is autocorrelation at lag k , κ is integrated autocorrelation time).

Inference using Posterior Samples from MCMC Runs

A powerful feature of MCMC and the Bayesian approach is that all inference is based on the (samples from the) joint posterior distribution.

- We can therefore address a wide range of substantive questions by calculating appropriate summaries from the posterior samples.

Inference using Posterior Samples from MCMC Runs

A powerful feature of MCMC and the Bayesian approach is that all inference is based on the (samples from the) joint posterior distribution.

- We can therefore address a wide range of substantive questions by calculating appropriate summaries from the posterior samples.
- **Typical Reporting:**
 - Report either the **mean** or **median** of the posterior samples for each parameter of interest as a point estimate.
 - Report **2.5% and 97.5% percentiles** (quantiles) of the posterior sample for each parameter. These define a **95% posterior credible interval**.

Inference using Posterior Samples from MCMC Runs

A powerful feature of MCMC and the Bayesian approach is that all inference is based on the (samples from the) joint posterior distribution.

- We can therefore address a wide range of substantive questions by calculating appropriate summaries from the posterior samples.
- **Typical Reporting:**
 - Report either the **mean** or **median** of the posterior samples for each parameter of interest as a point estimate.
 - Report **2.5% and 97.5% percentiles** (quantiles) of the posterior sample for each parameter. These define a **95% posterior credible interval**.
- **Interpretation of 95% Credible Interval:**
 - An interval within which the true parameter value lies with 95% probability (given the model and data).

Derived Quantities from Posterior Samples

Once we have a sample from the joint posterior distribution, we can answer many interesting questions simply by performing calculations on these samples.

Derived Quantities from Posterior Samples

Once we have a sample from the joint posterior distribution, we can answer many interesting questions simply by performing calculations on these samples. **Examples:** For parameters $\theta, \theta_1, \theta_2$ (for which we have posterior samples):

- **What is the probability that $\theta > 0$?**
 - Calculate the proportion of posterior samples for θ that are greater than 0.

Derived Quantities from Posterior Samples

Once we have a sample from the joint posterior distribution, we can answer many interesting questions simply by performing calculations on these samples. **Examples:** For parameters $\theta, \theta_1, \theta_2$ (for which we have posterior samples):

- **What is the probability that $\theta > 0$?**
 - Calculate the proportion of posterior samples for θ that are greater than 0.
- **What is the probability that $\theta_1 > \theta_2$?**
 - Calculate the proportion of posterior samples where the sample for θ_1 is greater than the sample for θ_2 .

Derived Quantities from Posterior Samples

Once we have a sample from the joint posterior distribution, we can answer many interesting questions simply by performing calculations on these samples. **Examples:** For parameters $\theta, \theta_1, \theta_2$ (for which we have posterior samples):

- **What is the probability that $\theta > 0$?**
 - Calculate the proportion of posterior samples for θ that are greater than 0.
- **What is the probability that $\theta_1 > \theta_2$?**
 - Calculate the proportion of posterior samples where the sample for θ_1 is greater than the sample for θ_2 .

- **What is a 95% interval for a function of parameters, e.g., $\theta_1/(\theta_1 + \theta_2)$?**
 - For each set of posterior samples $(\theta_1^{(s)}, \theta_2^{(s)})$, calculate the derived quantity $d^{(s)} = \theta_1^{(s)}/(\theta_1^{(s)} + \theta_2^{(s)})$.
 - Then find the 2.5% and 97.5% quantiles of the resulting samples $\{d^{(s)}\}$.

This flexibility is a major strength of Bayesian inference with MCMC.

Model Comparison in MCMC

Recap: Frequentist Model Comparison Options

- Likelihood Ratio (Deviance) Tests
- Wald Tests
- Information Criteria – e.g., AIC (Akaike Information Criterion) / BIC (Bayesian Information Criterion)

Model Comparison in MCMC

Recap: Frequentist Model Comparison Options

- Likelihood Ratio (Deviance) Tests
- Wald Tests
- Information Criteria – e.g., AIC (Akaike Information Criterion) / BIC (Bayesian Information Criterion)

Bayesian Model Comparison with MCMC Output:

- Here we look at a criterion that can be used with MCMC output.
- For a linear regression model, one such criterion, the **Deviance Information Criterion (DIC)**, is roughly equivalent to AIC.
- *(Note: Modern Bayesian practice often favors PSIS-LOO or WAIC over DIC due to some theoretical and practical issues with DIC, but DIC is historically important and still seen.)*

Deviance Information Criterion (DIC): Concept

- A natural way to compare models is to use a criterion based on a **trade-off** between:
 - 1 The **fit** of the data to the model.
 - 2 The corresponding **complexity** of the model.

Deviance Information Criterion (DIC): Concept

- A natural way to compare models is to use a criterion based on a **trade-off** between:
 - 1 The **fit** of the data to the model.
 - 2 The corresponding **complexity** of the model.
- DIC does this in a Bayesian way:

DIC = 'goodness of fit' + 'complexity penalty'

Deviance Information Criterion (DIC): Concept

- A natural way to compare models is to use a criterion based on a **trade-off** between:
 - 1 The **fit** of the data to the model.
 - 2 The corresponding **complexity** of the model.
- DIC does this in a Bayesian way:

DIC = 'goodness of fit' + 'complexity penalty'

- **Fit is measured by deviance** $D(\theta)$:

$$D(\theta) = -2 \log L(\text{data}|\theta)$$

(Lower deviance means better fit).

Deviance Information Criterion (DIC): Concept

- A natural way to compare models is to use a criterion based on a **trade-off** between:
 - 1 The **fit** of the data to the model.
 - 2 The corresponding **complexity** of the model.
- DIC does this in a Bayesian way:

DIC = 'goodness of fit' + 'complexity penalty'

- **Fit is measured by deviance** $D(\theta)$:

$$D(\theta) = -2 \log L(\text{data}|\theta)$$

(Lower deviance means better fit).

DIC: Effective Number of Parameters (p_D)

Complexity is measured by p_D , defined as:

$$p_D = E_{\theta|\text{data}}[D(\theta)] - D(E_{\theta|\text{data}}[\theta])$$

Which can be written as:

$$p_D = \overline{D(\theta)} - D(\bar{\theta})$$

DIC: Effective Number of Parameters (p_D)

Complexity is measured by p_D , defined as:

$$p_D = E_{\theta|\text{data}}[D(\theta)] - D(E_{\theta|\text{data}}[\theta])$$

Which can be written as:

$$p_D = \overline{D(\theta)} - D(\bar{\theta})$$

Interpretation:

- $\overline{D(\theta)}$: The posterior mean deviance (average deviance over all posterior samples of θ).
- $D(\bar{\theta})$: The deviance evaluated at the posterior mean of the parameters ($E_{\theta|\text{data}}[\theta]$).

Interpretation:

- So, p_D is the "Posterior mean deviance minus the deviance evaluated at the posterior mean of the parameters."
- p_D acts as a penalty for model complexity. More flexible models that can fit noise will tend to have a larger difference between $\overline{D(\theta)}$ and $D(\bar{\theta})$, thus a larger p_D .

DIC (Continued): Definition

The DIC is then defined analogously to AIC as:

$$\text{DIC} = D(\bar{\theta}) + 2p_D$$

DIC (Continued): Definition

The DIC is then defined analogously to AIC as:

$$\text{DIC} = D(\bar{\theta}) + 2p_D$$

Alternatively, it can also be expressed using the posterior mean deviance:

$$\text{DIC} = \overline{D(\theta)} + p_D$$

(Since $p_D = \overline{D(\theta)} - D(\bar{\theta})$, substituting gives $D(\bar{\theta}) + 2(\overline{D(\theta)} - D(\bar{\theta})) = 2\overline{D(\theta)} - D(\bar{\theta})$. The form $\overline{D(\theta)} + p_D$ is more common for calculation from MCMC output.)

DIC (Continued): Definition

The DIC is then defined analogously to AIC as:

$$\text{DIC} = D(\bar{\theta}) + 2p_D$$

Alternatively, it can also be expressed using the posterior mean deviance:

$$\text{DIC} = \overline{D(\theta)} + p_D$$

(Since $p_D = \overline{D(\theta)} - D(\bar{\theta})$, substituting gives $D(\bar{\theta}) + 2(\overline{D(\theta)} - D(\bar{\theta})) = 2\overline{D(\theta)} - D(\bar{\theta})$. The form $\overline{D(\theta)} + p_D$ is more common for calculation from MCMC output.) **Interpretation:**

- Models with **smaller DIC** are better supported by the data (better balance of fit and complexity).

DIC (Continued): Definition

The DIC is then defined analogously to AIC as:

$$\text{DIC} = D(\bar{\theta}) + 2p_D$$

Alternatively, it can also be expressed using the posterior mean deviance:

$$\text{DIC} = \overline{D(\theta)} + p_D$$

(Since $p_D = \overline{D(\theta)} - D(\bar{\theta})$, substituting gives $D(\bar{\theta}) + 2(\overline{D(\theta)} - D(\bar{\theta})) = 2\overline{D(\theta)} - D(\bar{\theta})$. The form $\overline{D(\theta)} + p_D$ is more common for calculation from MCMC output.) **Interpretation:**

- Models with **smaller DIC** are better supported by the data (better balance of fit and complexity).

Deviance Information Criterion (DIC): Summary

- Diagnostic for model comparison.
- Goodness of fit criterion that is penalized for model complexity.
- Generalization of the Akaike Information Criterion (AIC), particularly where the effective number of parameters (df) is not clearly known beforehand (as in complex Bayesian models).
- Used for comparing non-nested models (e.g., models with the same number of parameters but different variables).
- Can be valuable for testing improved goodness of fit of non-linear models (e.g., Logit) where the likelihood (and hence deviance) might be defined differently or approximated.

Deviance Information Criterion (DIC): Summary

- Estimated by MCMC sampling; output typically includes:

Deviance Information Criterion (DIC): Summary

- Estimated by MCMC sampling; output typically includes:

Example Output Components:

- 'Dbar' ($\overline{D(\theta)}$): The average deviance from the complete set of (post-burn-in) MCMC iterations.
- 'D(thetabar)' ($D(\bar{\theta})$): The deviance at the expected value (posterior mean) of the unknown parameters.
- 'pD': The estimated effective degrees of freedom consumed in the fit (i.e., 'Dbar - D(thetabar)').
- 'DIC': Fit + Complexity (i.e., 'Dbar + pD' or 'D(thetabar) + 2*pD').

Rule of Thumb

NB: Lower DIC values = better, more parsimonious model.

Deviance Information Criterion (DIC): Summary

- **Caveat:** DIC is somewhat controversial and has known limitations (e.g., sensitivity to parameterization, issues with non-Normal posteriors).

Reference: Spiegelhalter, D.J., Best, N.G., Carlin, B.P. and van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society, Series B* 64: 583-640.

Some Guidance on Interpreting DIC Differences

- Any **decrease** in DIC suggests a better model.

Some Guidance on Interpreting DIC Differences

- Any **decrease** in DIC suggests a better model.
- **Stochastic Nature of MCMC:** With small differences in DIC, you should confirm if this is a real difference by checking results with:
 - Different MCMC random seeds.
 - Different starting values for the chains.

Some Guidance on Interpreting DIC Differences

- Any **decrease** in DIC suggests a better model.
- **Stochastic Nature of MCMC:** With small differences in DIC, you should confirm if this is a real difference by checking results with:
 - Different MCMC random seeds.
 - Different starting values for the chains.
- **Rules of Thumb (often borrowed from AIC experience):**
 - A model with a ΔDIC (difference from the best model's DIC) within **1-2** of the best model has substantial support in the data and should be considered along with the best model.
 - A ΔDIC value within **4-7** of the best model has considerably less support.
 - A ΔDIC value > 10 indicates that the worse model has virtually no support and can often be omitted from further consideration.
- *(These are heuristics and context always matters!)*

Outline for this Section

- 1 Introduction and Session Aims (9:30 – 10:00)
- 2 Theoretical – Bayes (Part 1) (10:00 – 11:30)
 - What is Bayesian Statistics?
 - Statistical Distributions
 - Priors
 - Example: Normal Data, Unknown Mean (Known Variance)
- 3 Theoretical – Bayes (Part 2) (11:50am – 12:10pm)
 - The Computational Challenge Intro to MCMC
- 4 Summary of Morning Session (12:10 – 12:30pm)
- 5 Practical Application (1:30 – 4:15pm)
 - Tools (Stan Overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison
- 6 Review and Closing (4:15 – 4:30pm)

Morning Session Recap

Covered foundational Bayesian concepts:

- Bayes' Theorem: updating beliefs with data.
- Priors, Likelihood, Posterior. Conjugacy as a special case.
- Credible Intervals for parameter uncertainty.
- Computational challenge leading to MCMC (e.g., Gibbs Sampling).

After lunch: Stan, MCMC diagnostics, practical regression!

Outline for this Section

- 1 Introduction and Session Aims (9:30 – 10:00)
- 2 Theoretical – Bayes (Part 1) (10:00 – 11:30)
 - What is Bayesian Statistics?
 - Statistical Distributions
 - Priors
 - Example: Normal Data, Unknown Mean (Known Variance)
- 3 Theoretical – Bayes (Part 2) (11:50am – 12:10pm)
 - The Computational Challenge Intro to MCMC
- 4 Summary of Morning Session (12:10 – 12:30pm)
- 5 Practical Application (1:30 – 4:15pm)**
 - Tools (Stan Overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison
- 6 Review and Closing (4:15 – 4:30pm)

Probabilistic Programming Languages (PPLs)

Software for specifying Bayesian models and performing inference.

- **Stan** (BSD-3 License) - Our focus.
- Turing.jl (MIT License)
- PyMC (Apache License)
- JAGS, BUGS (GPL License)



- High-performance platform using HMC/NUTS.
- Stan language (C++ like).
- Interfaces: R (rstan), Python, Julia, etc.

Stan Code Example: Linear Regression Blocks

```
1 data { /* Observed data and fixed quantities */ }
2 parameters { /* Parameters to estimate */ }
3 transformed parameters { /* Optional: Derived parameters */
  }
4 model { /* Priors and likelihood */ }
5 generated quantities { /* Optional: Predictions, log-
  likelihood */ }
```

Listing 4: Stan model structure.

Recommended References

- Gelman et al. (2013) - Ch 6: Model checking
- McElreath (2020) - Ch 4: Geocentric Models
- Gelman, Hill Vehtari (2020): Ch 6, Ch 11
- Gelman et al. (2020) - "Workflow Paper"

The Reality of Modeling

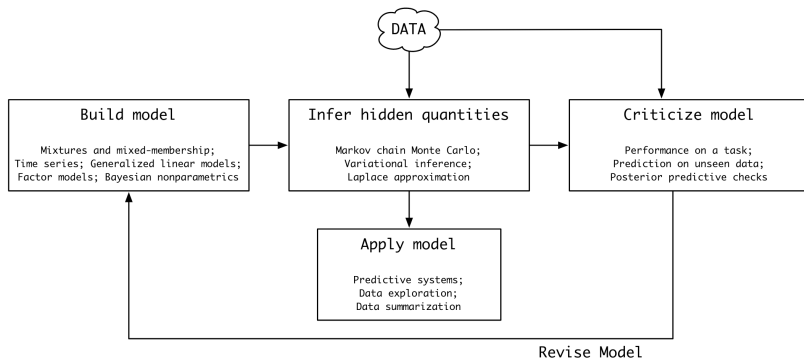


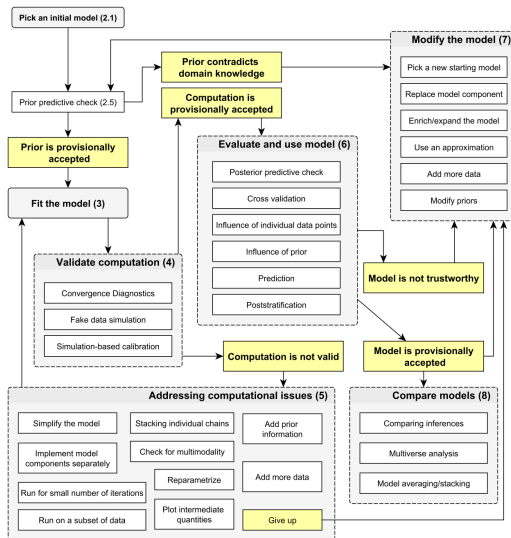
Figure: Follow a workflow!

Bayesian Workflow Overview

Iterative process (Gelman et al. 2020):

- 1 Understand problem → Formulate model
- 2 Implement model (Stan)
- 3 Prior predictive checks
- 4 Fit model (Run MCMC)
- 5 Assess MCMC convergence diagnostics
- 6 Posterior predictive checks
- 7 Evaluate, critique, compare, improve → Iterate

Bayesian Workflow Overview



MCMC Diagnostics: Did it Work?

Check if chains have: Converged, Mixed well, Produced enough effective samples.

- **Burn-in (Warmup):** Discard initial iterations.
- **Trace Plots:** Parameter value vs. iteration. Look for "fuzzy caterpillar" (stationarity).
- **Posterior Density Plots:** Smoothed histogram of samples.

Key MCMC Diagnostics (Continued)

- **R-hat (\hat{R}):** Potential scale reduction factor. $\hat{R} \approx 1$ (e.g., ≤ 1.01) for convergence (requires multiple chains).
- **Autocorrelation (ACF):** Correlation between samples at different lags. Want it to drop quickly.
- **Effective Sample Size (ESS / N_{eff}):** Number of independent-equivalent samples. Want ESS high (e.g., ≥ 100 -400 per chain for reliable estimates).
- **Monte Carlo Standard Error (MCSE):** MCMC sampling error for posterior summaries. Want it small.

Recommended References: Linear Regression

Recommended References

- Gelman et al. (2013): Ch 14, Ch 16
- McElreath (2020) - Ch 4
- Gelman, Hill Vehtari (2020): Ch 7, Ch 8, Ch 10

Linear Regression Everywhere



Bayesian Linear Regression: Model Structure

We specify a Bayesian linear regression by defining its components:

1. Likelihood (Data Generating Process): Describes how we believe the observed outcome y_i is generated, given the model parameters. For linear regression, we typically assume a Normal (Gaussian) distribution for the errors:

$$y_i \sim \text{Normal}(\mu_i, \sigma)$$

- y_i : The outcome for observation i .
- μ_i : The expected value (mean) of y_i .
- σ : The standard deviation of the errors (residuals), assumed constant across observations (homoscedasticity).

Bayesian Linear Regression: Model Structure

We specify a Bayesian linear regression by defining its components:

1. Likelihood (Data Generating Process): Describes how we believe the observed outcome y_i is generated, given the model parameters. For linear regression, we typically assume a Normal (Gaussian) distribution for the errors:

$$y_i \sim \text{Normal}(\mu_i, \sigma)$$

- y_i : The outcome for observation i .
- μ_i : The expected value (mean) of y_i .
- σ : The standard deviation of the errors (residuals), assumed constant across observations (homoscedasticity).

2. Linear Predictor (μ_i): The expected value μ_i is modeled as a linear combination of predictors:

$$\mu_i = \alpha + \mathbf{X}_i\boldsymbol{\beta}$$

- α : The intercept term (expected y_i when all predictors in \mathbf{X}_i are zero).
- \mathbf{X}_i : A row vector of predictor values for observation i .
- $\boldsymbol{\beta}$: A column vector of regression coefficients (slopes), indicating the change in μ_i for a one-unit change in the corresponding predictor in \mathbf{X}_i .

Bayesian Linear Regression: Priors for Parameters

In the Bayesian framework, every unknown parameter requires a **prior distribution**, reflecting our beliefs about the parameter *before* seeing the data.

3. Priors for Model Parameters:

- **Intercept (α):**

$$\alpha \sim \text{Prior for intercept}$$

A common choice is a weakly informative Normal prior.

Bayesian Linear Regression: Priors for Parameters

In the Bayesian framework, every unknown parameter requires a **prior distribution**, reflecting our beliefs about the parameter *before* seeing the data.

3. Priors for Model Parameters:

- **Intercept (α):**

$$\alpha \sim \text{Prior for intercept}$$

A common choice is a weakly informative Normal prior.

- **Coefficients (β):**

$$\beta \sim \text{Prior for coefficients}$$

Often, each β_k is given an independent, weakly informative Normal prior, e.g., $N(0, 2.5)$ or $N(0, 5)$, especially if predictors are standardized.

Bayesian Linear Regression: Priors for Parameters

In the Bayesian framework, every unknown parameter requires a **prior distribution**, reflecting our beliefs about the parameter *before* seeing the data.

3. Priors for Model Parameters:

- **Intercept (α):**

$$\alpha \sim \text{Prior for intercept}$$

A common choice is a weakly informative Normal prior.

- **Coefficients (β):**

$$\beta \sim \text{Prior for coefficients}$$

Often, each β_k is given an independent, weakly informative Normal prior, e.g., $N(0, 2.5)$ or $N(0, 5)$, especially if predictors are standardized.

Bayesian Linear Regression: Priors for Parameters

In the Bayesian framework, every unknown parameter requires a **prior distribution**, reflecting our beliefs about the parameter *before* seeing the data.

3. Priors for Model Parameters:

- **Intercept (α):**

$$\alpha \sim \text{Prior for intercept}$$

A common choice is a weakly informative Normal prior.

- **Coefficients (β):**

$$\beta \sim \text{Prior for coefficients}$$

Often, each β_k is given an independent, weakly informative Normal prior, e.g., $N(0, 2.5)$ or $N(0, 5)$, especially if predictors are standardized.

3. Priors for Model Parameters:

- **Error Scale (σ):**

$\sigma \sim$ Prior for error scale (must be positive)

Since σ must be positive, common choices include Half-Normal, Half-Cauchy, Half-Student-t, or Exponential priors, e.g., Half-Cauchy(0, 5).

3. Priors for Model Parameters:

- **Error Scale (σ):**

$\sigma \sim$ Prior for error scale (must be positive)

Since σ must be positive, common choices include Half-Normal, Half-Cauchy, Half-Student-t, or Exponential priors, e.g., Half-Cauchy(0, 5).

Combining Components

The likelihood and priors are combined via Bayes' theorem to obtain the posterior distribution $P(\alpha, \beta, \sigma | y, \mathbf{X})$, from which we make inferences.

Bayesian Linear Regression Specification

$$y_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \mathbf{X}_i\boldsymbol{\beta}$$

$\alpha \sim$ Prior for intercept

$\boldsymbol{\beta} \sim$ Prior for coefficients

$\sigma \sim$ Prior for error scale (positive)

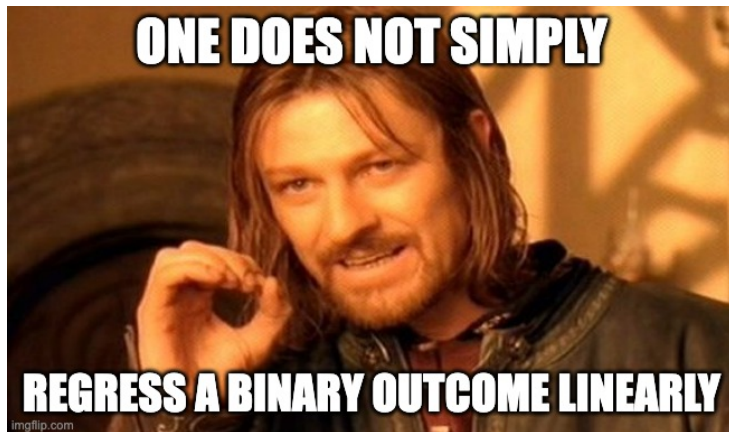
Stan Code: Linear Regression (Practical)

```
1 data {
2   int<lower=0> N; int<lower=0> K;
3   matrix[N, K] X; vector[N] y;
4 }
5 parameters {
6   real alpha; vector[K] beta; real<lower=0> sigma;
7 }
8 model {
9   alpha ~ normal(0, 10);
10  beta ~ normal(0, 2.5);
11  sigma ~ cauchy(0, 5);
12  y ~ normal(X * beta + alpha, sigma);
13 }
14 generated quantities {
15   vector[N] y_rep; vector[N] log_lik;
16   for (n in 1:N) {
17     real mu_n = X[n] * beta + alpha;
18     y_rep[n] = normal_rng(mu_n, sigma);
19     log_lik[n] = normal_lpdf(y[n] | mu_n, sigma);
20   }
21 }
```

Recommended References: Logistic Regression

Recommended References

- Gelman et al. (2013) - Ch 16: GLMs
- McElreath (2020) - Ch 10 (GLM), Ch 11 (Binomial)
- Gelman, Hill Vehtari (2020): Ch 13-15



Bayesian Logistic Regression: Model Structure (Binary Outcome)

When the outcome y_i is binary (0 or 1, e.g., success/failure, yes/no), we use logistic regression.

1. Likelihood (Data Generating Process for Binary Data): Each observation y_i is assumed to follow a **Bernoulli distribution**, governed by a probability p_i :

$$y_i \sim \text{Bernoulli}(p_i)$$

- y_i : The binary outcome for observation i (either 0 or 1).
- p_i : The probability that $y_i = 1$ (e.g., probability of success). This p_i must be between 0 and 1.

Bayesian Logistic Regression: Model Structure (Binary Outcome)

When the outcome y_i is binary (0 or 1, e.g., success/failure, yes/no), we use logistic regression.

1. Likelihood (Data Generating Process for Binary Data): Each observation y_i is assumed to follow a **Bernoulli distribution**, governed by a probability p_i :

$$y_i \sim \text{Bernoulli}(p_i)$$

- y_i : The binary outcome for observation i (either 0 or 1).
- p_i : The probability that $y_i = 1$ (e.g., probability of success). This p_i must be between 0 and 1.

Bayesian Logistic Regression: Model Structure (Binary Outcome)

2. Link Function (Connecting Predictors to Probability p_i): We can't directly model $p_i = \alpha + \mathbf{X}_i\beta$ because the linear predictor can range from $-\infty$ to $+\infty$, while p_i is constrained to $[0, 1]$.

- We use a **link function** to map the linear predictor to the probability scale.
- For logistic regression, this is the **logit link**:

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = \eta_i$$

- η_i is the linear predictor (on the log-odds scale).

Bayesian Logistic Regression: Model Structure (Binary Outcome)

2. Link Function (Connecting Predictors to Probability p_i): We can't directly model $p_i = \alpha + \mathbf{X}_i\beta$ because the linear predictor can range from $-\infty$ to $+\infty$, while p_i is constrained to $[0, 1]$.

- We use a **link function** to map the linear predictor to the probability scale.
- For logistic regression, this is the **logit link**:

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = \eta_i$$

- η_i is the linear predictor (on the log-odds scale).

Bayesian Logistic Regression: Model Structure (Binary Outcome)

3. Linear Predictor (η_i on Log-Odds Scale): The log-odds of success, η_i , is modeled as a linear combination of predictors:

$$\eta_i = \alpha + \mathbf{X}_i\beta$$

- α : Intercept on the log-odds scale.
- \mathbf{X}_i : Predictor values for observation i .
- β : Coefficients (slopes) on the log-odds scale.

Bayesian Logistic Regression: Priors Interpretation

Similar to linear regression, we need priors for the unknown parameters α and β .

4. Priors for Model Parameters (on Log-Odds Scale):

- **Intercept (α):**

$\alpha \sim$ Prior for intercept (log-odds scale)

A common choice is a weakly informative Normal prior, e.g., $N(0, 5)$ or $N(0, 2.5)$. A wider prior like $N(0, 10)$ might be used if less is known about the baseline log-odds.

Bayesian Logistic Regression: Priors Interpretation

Similar to linear regression, we need priors for the unknown parameters α and β .

4. Priors for Model Parameters (on Log-Odds Scale):

- **Intercept (α):**

$\alpha \sim \text{Prior for intercept (log-odds scale)}$

A common choice is a weakly informative Normal prior, e.g., $N(0, 5)$ or $N(0, 2.5)$. A wider prior like $N(0, 10)$ might be used if less is known about the baseline log-odds.

- **Coefficients (β):**

$\beta \sim \text{Prior for coefficients (log-odds scale)}$

Often, each β_k is given an independent, weakly informative Normal prior, e.g., $N(0, 2.5)$. This suggests that a one-unit change in a (standardized) predictor is unlikely to change the log-odds by more than a few units.

Bayesian Logistic Regression: Priors Interpretation

Similar to linear regression, we need priors for the unknown parameters α and β .

4. Priors for Model Parameters (on Log-Odds Scale):

- **Intercept (α):**

$\alpha \sim \text{Prior for intercept (log-odds scale)}$

A common choice is a weakly informative Normal prior, e.g., $N(0, 5)$ or $N(0, 2.5)$. A wider prior like $N(0, 10)$ might be used if less is known about the baseline log-odds.

- **Coefficients (β):**

$\beta \sim \text{Prior for coefficients (log-odds scale)}$

Often, each β_k is given an independent, weakly informative Normal prior, e.g., $N(0, 2.5)$. This suggests that a one-unit change in a (standardized) predictor is unlikely to change the log-odds by more than a few units.

4. Priors for Model Parameters (on Log-Odds Scale): Interpretation of Coefficients:

- The parameters α and β_k are estimated on the **log-odds scale**.
 - β_k : The change in the log-odds of $y_i = 1$ for a one-unit increase in predictor X_{ik} , holding other predictors constant.

4. Priors for Model Parameters (on Log-Odds Scale): Interpretation of Coefficients:

- The parameters α and β_k are estimated on the **log-odds scale**.
 - β_k : The change in the log-odds of $y_i = 1$ for a one-unit increase in predictor X_{ik} , holding other predictors constant.
- For more intuitive interpretation, we often **exponentiate** the coefficients to get **Odds Ratios (ORs)**:
 - e^{β_k} : The multiplicative change in the *odds* of $y_i = 1$ for a one-unit increase in X_{ik} .
 - $OR = 1$: No effect.
 - $OR > 1$: Increased odds.
 - $OR < 1$: Decreased odds.

Full Model

$y_i \sim \text{Bernoulli}(p_i)$, where $p_i = \text{invlogit}(\alpha + \mathbf{X}_i\beta)$. Priors are placed on α and β .

Bayesian Logistic Regression Specification

$$y_i \sim \text{Bernoulli}(p_i)$$

$$\text{logit}(p_i) = \eta_i = \alpha + \mathbf{X}_i\boldsymbol{\beta}$$

$\alpha \sim$ Prior for intercept (log-odds scale)

$\boldsymbol{\beta} \sim$ Prior for coefficients (log-odds scale)

Interpret coefficients as log-odds or exponentiate for Odds Ratios.

Stan Code: Logistic Regression (Practical)

```
1 data {
2   int<lower=0> N; int<lower=0> K;
3   matrix[N, K] X; int<lower=0, upper=1> y[N];
4 }
5 parameters {
6   real alpha; vector[K] beta;
7 }
8 model {
9   alpha ~ normal(0, 5); // Prior on log-odds
10  beta ~ normal(0, 2.5); // Prior on log-odds
11  y ~ bernoulli_logit(X * beta + alpha);
12 }
13 generated quantities {
14   vector[N] y_rep; vector[N] log_lik;
15   for (n in 1:N) {
16     real eta_n = X[n] * beta + alpha;
17     y_rep[n] = bernoulli_logit_rng(eta_n);
18     log_lik[n] = bernoulli_logit_lpmf(y[n] | eta_n);
19   }
20 }
```


Recommended References

- Gelman et al. (2013) - Ch 7
- Gelman, Hill Vehtari (2020) - Ch 11.8
- McElreath (2020) - Ch 7.5
- Vehtari, Gelman Gabry (2017) - LOO paper
- Spiegelhalter et al. (2002) - DIC paper
- Watanabe Opper (2010) - WAIC paper

The Pitfall of Simple Metrics



Figure: Model selection needs careful thought.

Model Comparison Overview

Goal: Estimate out-of-sample predictive accuracy.

- **LOO-CV (PSIS-LOO)**: Generally preferred. Check Pareto \hat{k} .
- **WAIC**: Asymptotically similar to LOO.
- **DIC**: Historically used. **Lower DIC is better.**

Calculated from *log-lik* in *generated quantities*.

Deviance Information Criterion (DIC)

DIC = "goodness of fit" + "complexity"

- Fit: $D(\bar{\theta}) = -2 \log L(\text{data}|\bar{\theta})$
- Complexity: $p_D = \overline{D(\theta)} - D(\bar{\theta})$
- DIC = $\overline{D(\theta)} + p_D$.

Using DIC: House Price Example (Bristol)

Comparing Models with DIC

Model Description	p_D	DIC
0: Null one-level	2.00	10728.31
1: Random intercepts	44.43	10498.74
3: Random intercepts and slopes	65.24	9807.47

Model 3 (random intercepts slopes) has lowest DIC. $\Delta DIC > 10$ often considered strong evidence.

LOO-CV and WAIC with Stan

Use R package **loo** with Stan fit object.

- Compare *elpd-loo* or *elpd-waic* (higher is better).
- *loo-compare()* for comparing multiple models.

Outline for this Section

- 1 Introduction and Session Aims (9:30 – 10:00)
- 2 Theoretical – Bayes (Part 1) (10:00 – 11:30)
 - What is Bayesian Statistics?
 - Statistical Distributions
 - Priors
 - Example: Normal Data, Unknown Mean (Known Variance)
- 3 Theoretical – Bayes (Part 2) (11:50am – 12:10pm)
 - The Computational Challenge Intro to MCMC
- 4 Summary of Morning Session (12:10 – 12:30pm)
- 5 Practical Application (1:30 – 4:15pm)
 - Tools (Stan Overview)
 - Bayesian Workflow
 - MCMC Diagnostics in Practice
 - Linear Regression with Stan
 - Logistic Regression with Stan
 - Model Comparison
- 6 Review and Closing (4:15 – 4:30pm)

Our Journey Through Bayesian Modeling:

- We started with the foundational concepts: What is probability from a Bayesian view? How does Bayes' Theorem allow us to update beliefs?
- We explored key components: priors (our initial knowledge), likelihoods (information from data), and the resulting posteriors (our updated knowledge).
- We saw the computational engine: MCMC (and specifically HMC/NUTS in Stan) enables us to tackle complex, realistic models where analytical solutions are out of reach.
- We put theory into practice with linear and logistic regression, learning how to specify models in Stan, check diagnostics, and interpret results.

Our Journey Through Bayesian Modeling:

- We started with the foundational concepts: What is probability from a Bayesian view? How does Bayes' Theorem allow us to update beliefs?
- We explored key components: priors (our initial knowledge), likelihoods (information from data), and the resulting posteriors (our updated knowledge).
- We saw the computational engine: MCMC (and specifically HMC/NUTS in Stan) enables us to tackle complex, realistic models where analytical solutions are out of reach.
- We put theory into practice with linear and logistic regression, learning how to specify models in Stan, check diagnostics, and interpret results.

Key Takeaways for Your Future Work:

- Bayesian statistics offers a powerful and flexible framework for nuanced data analysis and uncertainty quantification.
- The "Bayesian workflow" is an iterative process of model building, checking, and refinement – embrace it!
- Priors matter: Be thoughtful and transparent about your choices. Weakly informative priors are often a good start.
- Tools like Stan make advanced Bayesian methods accessible.

Course Reflection Moving Forward

Key Takeaways for Your Future Work:

- Bayesian statistics offers a powerful and flexible framework for nuanced data analysis and uncertainty quantification.
- The "Bayesian workflow" is an iterative process of model building, checking, and refinement – embrace it!
- Priors matter: Be thoughtful and transparent about your choices. Weakly informative priors are often a good start.
- Tools like Stan make advanced Bayesian methods accessible.

The Journey Continues!

This course is an introduction. The world of Bayesian modeling is vast and deep. Keep learning, keep practicing, and don't hesitate to explore more advanced topics and models.

Further Learning & Resources

- **Books:** Gelman et al. (BDA3), McElreath (Statistical Rethinking), Gelman, Hill, Vehtari (ROS).
- **Stan Resources:** User's Guide, Forums (<https://discourse.mc-stan.org/>).
- R packages: 'rstan', 'loo', 'bayesplot'.

“Today’s posterior is tomorrow’s prior.”

— Dennis V. Lindley

Questions & Thank You!

Thank you for your engagement and curiosity!
`diego.perezruiz@manchester.ac.uk`

Bibliography I



Bertsekas, D.P. and Tsitsiklis, J.N. (2008) *Introduction to Probability, 2nd Edition*. Athena Scientific.



Betancourt, M. (2019) *Probabilistic Building Blocks*. Available at:
https://betanalpha.github.io/assets/case_studies/probability_densities.html.



Blei, D.M. (2014) "Build, Compute, Critique, Repeat: Data Analysis with Latent Variable Models," *Annual Review of Statistics and Its Application*, 1(1), pp. 203–232.



Box, G.E.P. (1976) "Science and Statistics," *Journal of the American Statistical Association*, 71(356), pp. 791–799.



Carpenter, B. et al. (2017) "Stan : A Probabilistic Programming Language," *Journal of Statistical Software*, 76(1).



Dekking, F.M. et al. (2010) *A Modern Introduction to Probability and Statistics: Understanding Why and How*. Springer.



Diaconis, P. and Skyrms, B. (2019) *Ten Great Ideas about Chance*. Princeton University Press.



Ge, H., Xu, K. and Ghahramani, Z. (2018) "Turing: A Language for Flexible Probabilistic Inference," in *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 1682–1690.



Geisser, S. and Eddy, W.F. (1979) "A predictive approach to model selection," *Journal of the American Statistical Association*, 74(365), pp. 153–160.

Bibliography II



Gelfand, A.E., Dey, D.K. and Chang, H. (1992) "Model determination using predictive distributions with implementation via sampling-based methods," *Bayesian Statistics*. Edited by J.M. Bernardo et al. Oxford University Press.



Gelfand, A.E. (1996) "Model determination using sampling-based methods," *Markov chain Monte Carlo in practice*, pp. 145–161.



Gelman, A. et al. (2013) *Bayesian Data Analysis*. Chapman and Hall/CRC. (BDA3)



Gelman, A., Hill, J. and Vehtari, A. (2020) *Regression and Other Stories*. Cambridge University Press. (ROS)



Gelman, A. et al. (2020) *Bayesian Workflow*. Available at: <http://arxiv.org/abs/2011.01808>.



Grimmett, G. and Stirzaker, D. (2020) *Probability and Random Processes: Fourth Edition*. Oxford University Press.



Jaynes, E.T. (2003) *Probability Theory: The Logic of Science*. Cambridge university press.



Khan, M.E. and Rue, H. (2021) *The Bayesian Learning Rule*. Available at: <http://arxiv.org/abs/2107.04562>.



Kurt, W. (2019) *Bayesian Statistics the Fun Way: Understanding Statistics and Probability with Star Wars, LEGO, and Rubber Ducks*. No Starch Press.



Van Der Linde, A. (2005) "DIC in variable selection," *Statistica Neerlandica*, 59(1), pp. 45–56.

Bibliography III



McElreath, R. (2020) *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. CRC press.



Salvatier, J., Wiecki, T.V. and Fonnesbeck, C. (2016) "Probabilistic programming in Python using PyMC3," *PeerJ Computer Science*, 2, p. e55.



Schoot, R. van de et al. (2021) "Bayesian Statistics and Modelling," *Nature Reviews Methods Primers*, 1(1), pp. 1–26.



Spiegelhalter, D.J., Best, N.G., Carlin, B.P. and van der Linde, A. (2002) "Bayesian measures of model complexity and fit," *Journal of the royal statistical society: Series B (statistical methodology)*, 64(4), pp. 583–639.



Vehtari, A., Gelman, A. and Gabry, J. (2017) "Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC." *Statistics and Computing*, 27(5), pp. 1413-1432.



Watanabe, S. and Opper, M. (2010) "Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory.," *Journal of machine learning research*, 11(Dec), pp. 3571-3594.