

# Mathematical Foundations and Detailed Derivations for the Analysis Pipeline

Automated Exposition

December 4, 2025

## Abstract

This document integrates and exhaustively explains every mathematical and analytical component used in the provided analysis code. The exposition covers: the Pearson correlation and permutation test; ordinary least squares and residual sums of squares; the Bayesian Information Criterion (BIC) and its Laplace approximation to Bayes factors; the Savage–Dickey density ratio and practical estimation of posterior density at a point using Gaussian kernel density estimation (KDE); Bayesian linear regression modeling (centered and non-centered parameterizations); Student- $t$  likelihood for robust regression; Markov chain Monte Carlo (MCMC) diagnostics including divergences, Gelman–Rubin  $\hat{R}$ , and effective sample size; posterior predictive checks (PPC); leave-one-out cross-validation (LOO) and Pareto-smoothed importance sampling (PSIS); and the conservative decision heuristics used to combine evidence. Each formula is derived or motivated and practical numerical considerations are discussed.

## Contents

<b>1 Notation and basic setup</b>	<b>4</b>
<b>2 Pearson correlation: definition and algebraic properties</b>	<b>4</b>
2.1 Definition . . . . .	4
2.2 Relation to linear regression . . . . .	4
2.3 Properties . . . . .	5
<b>3 Permutation test for Pearson correlation</b>	<b>5</b>
3.1 Null hypothesis and permutation principle . . . . .	5
3.2 Test statistic . . . . .	5
3.3 Permutation distribution . . . . .	5
3.4 Two-sided p-value with small-sample correction . . . . .	5
3.5 Remarks on interpretation . . . . .	5
<b>4 Ordinary least squares, residual sum of squares, and <math>R^2</math></b>	<b>6</b>
4.1 OLS estimators . . . . .	6
4.2 Residual sum of squares (RSS) . . . . .	6
4.3 Coefficient of determination $R^2$ . . . . .	6

<b>5 Bayesian Information Criterion (BIC) and Laplace approximation to Bayes factors</b>	<b>6</b>
5.1 BIC definition . . . . .	6
5.2 Laplace approximation and Bayes factor . . . . .	7
5.3 Interpretation of BF values . . . . .	7
<b>6 Savage–Dickey density ratio: exact identity and practical estimation</b>	<b>8</b>
6.1 Savage–Dickey identity . . . . .	8
6.2 Prior density at zero for a normal prior . . . . .	8
6.3 Posterior density at zero: estimation from samples . . . . .	8
6.4 Bandwidth selection and practical issues . . . . .	8
6.5 Interpretation and limitations . . . . .	9
<b>7 Bayesian linear regression: model specification and parameterizations</b>	<b>9</b>
7.1 Standard (centered) parameterization . . . . .	9
7.2 Non-centered parameterization . . . . .	9
7.3 Why non-centered helps . . . . .	9
<b>8 Student–<i>t</i> likelihood for robust regression</b>	<b>9</b>
8.1 Student– <i>t</i> density . . . . .	9
8.2 Robustness properties . . . . .	10
8.3 Priors for $\nu$ . . . . .	10
<b>9 MCMC sampling diagnostics</b>	<b>10</b>
9.1 Divergences in Hamiltonian Monte Carlo (NUTS) . . . . .	10
9.2 Gelman–Rubin $\hat{R}$ statistic . . . . .	10
9.3 Effective sample size (ESS) . . . . .	11
<b>10 Posterior predictive checks (PPC)</b>	<b>11</b>
10.1 Posterior predictive distribution . . . . .	11
10.2 Generating replicated datasets . . . . .	11
10.3 Summary statistics and discrepancy measures . . . . .	12
<b>11 Leave-one-out cross-validation (LOO) and PSIS</b>	<b>12</b>
11.1 LOO predictive density . . . . .	12
11.2 Importance sampling identity . . . . .	12
11.3 Pareto-smoothed importance sampling (PSIS) . . . . .	12
11.4 LOO model comparison . . . . .	12
<b>12 Highest density interval (HDI) and credible intervals</b>	<b>13</b>
12.1 Definition . . . . .	13
12.2 HDI exclusion of zero . . . . .	13
<b>13 Combining evidence: heuristics used in the pipeline</b>	<b>13</b>
<b>14 Numerical and practical considerations</b>	<b>14</b>
14.1 Finite precision and floors . . . . .	14
14.2 Monte Carlo variability . . . . .	14
14.3 Standardization of predictors . . . . .	14

<b>15 Worked symbolic derivations and small proofs</b>	<b>14</b>
15.1 Derivation: relation between $r$ and OLS slope . . . . .	14
15.2 Derivation: BIC for Gaussian linear model . . . . .	14
15.3 Derivation: Savage–Dickey identity (sketch) . . . . .	15
<b>16 Summary of recommended practical workflow (algorithmic steps)</b>	<b>15</b>
<b>17 Appendix: useful formulae and constants</b>	<b>16</b>

# 1 Notation and basic setup

Let  $\{(x_i, y_i)\}_{i=1}^n$  denote paired observations. Vectors  $x = (x_1, \dots, x_n)^\top$  and  $y = (y_1, \dots, y_n)^\top$  are column vectors. We use the following standard sample statistics:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Define centered variables

$$\tilde{x}_i = x_i - \bar{x}, \quad \tilde{y}_i = y_i - \bar{y}.$$

Define sums of squares and cross-products

$$S_{xx} = \sum_{i=1}^n \tilde{x}_i^2, \quad S_{yy} = \sum_{i=1}^n \tilde{y}_i^2, \quad S_{xy} = \sum_{i=1}^n \tilde{x}_i \tilde{y}_i.$$

We denote by  $\|\cdot\|$  the Euclidean norm and by  $\mathbf{1}\{\cdot\}$  the indicator function.

# 2 Pearson correlation: definition and algebraic properties

## 2.1 Definition

The sample Pearson correlation coefficient  $r$  between  $x$  and  $y$  is

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}.$$

This is the sample estimate of the population Pearson correlation  $\rho$ .

## 2.2 Relation to linear regression

Consider the simple linear regression model with intercept:

$$y_i = \alpha + \beta x_i + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2).$$

The ordinary least squares (OLS) slope estimator (when regressing  $y$  on  $x$  with intercept) is

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}}.$$

The sample correlation relates to  $\hat{\beta}$  via

$$\hat{\beta} = r \cdot \frac{\sqrt{S_{yy}}}{\sqrt{S_{xx}}}.$$

Equivalently,

$$r = \hat{\beta} \cdot \frac{\sqrt{S_{xx}}}{\sqrt{S_{yy}}}.$$

## 2.3 Properties

- $r \in [-1, 1]$ .
- $r = 1$  or  $-1$  indicates perfect linear relationship.
- $r$  is invariant to separate location shifts of  $x$  and  $y$  and to separate positive scaling.

## 3 Permutation test for Pearson correlation

### 3.1 Null hypothesis and permutation principle

The permutation test assesses the null hypothesis  $H_0$ : “no association between  $x$  and  $y$ ” (exchangeability of  $y$  relative to  $x$ ). Under  $H_0$ , the joint distribution of  $(x, y)$  factorizes so that the labels of  $y$  can be permuted without changing the distribution of the data.

### 3.2 Test statistic

Use the observed Pearson correlation  $r_{\text{obs}}$  as the test statistic:

$$r_{\text{obs}} = r(x, y).$$

### 3.3 Permutation distribution

Generate  $N_{\text{perm}}$  random permutations  $\pi_j$  of indices  $\{1, \dots, n\}$ . For each permutation compute

$$r^{(j)} = r(x, y_{\pi_j}).$$

This yields an empirical null distribution of the statistic under  $H_0$ .

### 3.4 Two-sided p-value with small-sample correction

Count exceedances:

$$C = \sum_{j=1}^{N_{\text{perm}}} \mathbf{1}\{|r^{(j)}| \geq |r_{\text{obs}}|\}.$$

The permutation p-value with the standard Monte Carlo correction is

$$p = \frac{C + 1}{N_{\text{perm}} + 1}.$$

This correction ensures that  $p$  is never zero and yields an unbiased estimator of the permutation p-value in the Monte Carlo sense.

### 3.5 Remarks on interpretation

- The permutation test is nonparametric and exact (conditional on the observed  $x$ ) if all permutations are enumerated; with Monte Carlo sampling it is approximate.
- The test assesses association but not causation.

- The test is sensitive to monotone linear association as measured by Pearson  $r$ ; it is not optimal for nonlinear associations.

## 4 Ordinary least squares, residual sum of squares, and $R^2$

### 4.1 OLS estimators

In matrix form, for a model with intercept and single predictor  $x$ , define design matrix

$$X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \in \mathbb{R}^{n \times 2}.$$

The OLS estimator is

$$\hat{\beta}_{\text{vec}} = (X^\top X)^{-1} X^\top y = \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}.$$

### 4.2 Residual sum of squares (RSS)

Predicted values  $\hat{y} = X\hat{\beta}_{\text{vec}}$ . Residuals  $e = y - \hat{y}$ . Residual sum of squares:

$$\text{RSS} = \sum_{i=1}^n e_i^2 = \|y - X\hat{\beta}_{\text{vec}}\|^2.$$

### 4.3 Coefficient of determination $R^2$

Total sum of squares:

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2 = S_{yy}.$$

Explained sum of squares ESS = TSS – RSS. Then

$$R^2 = \frac{\text{ESS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}.$$

For simple linear regression,  $R^2 = r^2$ .

## 5 Bayesian Information Criterion (BIC) and Laplace approximation to Bayes factors

### 5.1 BIC definition

For a model  $M$  with  $k$  free parameters (counting intercepts and variance parameters as appropriate) and sample size  $n$ , the Bayesian Information Criterion (Schwarz, 1978) is

$$\text{BIC} = -2 \log \hat{L} + k \log n$$

where  $\hat{L}$  is the maximized likelihood. For Gaussian linear models with residual sum of squares RSS, the log-likelihood (up to additive constants not depending on parameters) is

$$\log L \propto -\frac{n}{2} \log \left( \frac{\text{RSS}}{n} \right).$$

Thus a commonly used expression for BIC in the Gaussian regression context is

$$\boxed{\text{BIC} = n \log \left( \frac{\text{RSS}}{n} \right) + k \log n.}$$

This is the form used in the code.

## 5.2 Laplace approximation and Bayes factor

The marginal likelihood (evidence) for model  $M$  is

$$p(y | M) = \int p(y | \theta, M) p(\theta | M) d\theta.$$

Laplace's method approximates this integral by expanding the log posterior around its mode and approximating by a Gaussian integral. The BIC arises as an asymptotic approximation to  $-2 \log p(y | M)$  up to constants. For two models  $M_0$  (null) and  $M_1$  (alternative), the Bayes factor in favor of  $M_0$  over  $M_1$  is

$$\text{BF}_{01} = \frac{p(y | M_0)}{p(y | M_1)}.$$

Using the BIC approximation,

$$-2 \log p(y | M) \approx \text{BIC} + \text{constant},$$

so

$$\log \text{BF}_{01} \approx -\frac{1}{2} \text{BIC}_0 + \frac{1}{2} \text{BIC}_1,$$

hence

$$\boxed{\text{BF}_{01} \approx \exp \left( \frac{\text{BIC}_0 - \text{BIC}_1}{2} \right).}$$

This is a large-sample approximation; it depends on the choice of parameterization and is less reliable for small  $n$  or weak priors.

## 5.3 Interpretation of BF values

Common interpretive thresholds (Jeffreys, Kass & Raftery) are:

- $\text{BF}_{01} > 10$ : strong evidence for  $M_0$ .
- $\text{BF}_{01} \in (3, 10]$ : moderate evidence for  $M_0$ .
- $\text{BF}_{01} \in (1/3, 3)$ : inconclusive.
- $\text{BF}_{01} < 1/10$ : strong evidence for  $M_1$ .

These are heuristic and context dependent.

## 6 Savage–Dickey density ratio: exact identity and practical estimation

### 6.1 Savage–Dickey identity

Consider nested models where  $M_0$  is obtained from  $M_1$  by fixing a parameter  $\beta$  to a point value (commonly 0). Suppose the prior under  $M_1$  factorizes as  $p(\beta, \phi) = p(\beta)p(\phi)$  and the prior under  $M_0$  is  $p(\phi)$  (i.e., the prior for nuisance parameters  $\phi$  is the same). Then the Bayes factor comparing  $M_0$  to  $M_1$  can be written as the ratio of prior to posterior marginal densities of  $\beta$  at the point:

$$\boxed{\text{BF}_{01} = \frac{p_{\text{prior}}(\beta = 0)}{p_{\text{post}}(\beta = 0)}}.$$

This is the Savage–Dickey density ratio. It is exact under the stated conditions.

### 6.2 Prior density at zero for a normal prior

If the prior for  $\beta$  is  $\beta \sim \mathcal{N}(0, \sigma_{\text{prior}}^2)$ , then the prior density at zero is

$$\boxed{p_{\text{prior}}(0) = \frac{1}{\sqrt{2\pi} \sigma_{\text{prior}}}}.$$

### 6.3 Posterior density at zero: estimation from samples

The posterior density  $p_{\text{post}}(\beta)$  is typically not available in closed form when using MCMC. Given posterior samples  $\{\beta^{(s)}\}_{s=1}^S$ , estimate  $p_{\text{post}}(0)$  using a kernel density estimator (KDE). For a Gaussian kernel  $K(u) = (2\pi)^{-1/2} \exp(-u^2/2)$  and bandwidth  $h > 0$ ,

$$\boxed{\hat{p}_{\text{post}}(0) = \frac{1}{Sh} \sum_{s=1}^S K\left(\frac{0 - \beta^{(s)}}{h}\right) = \frac{1}{Sh} \sum_{s=1}^S \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^{(s)})^2}{2h^2}\right)}.$$

Then the Savage–Dickey estimate is

$$\widehat{\text{BF}}_{01} = \frac{p_{\text{prior}}(0)}{\hat{p}_{\text{post}}(0)}.$$

### 6.4 Bandwidth selection and practical issues

- Bandwidth  $h$  controls bias–variance tradeoff. Common rules include Silverman’s rule of thumb:

$$h_{\text{Silverman}} = 0.9 \min\{\hat{\sigma}, \text{IQR}/1.34\} S^{-1/5},$$

where  $\hat{\sigma}$  is the sample standard deviation of  $\{\beta^{(s)}\}$  and IQR is the interquartile range.

- If posterior samples are few or concentrated, KDE can be unstable; a fallback is to approximate the posterior by a normal distribution with mean  $\hat{\mu}$  and standard deviation  $\hat{s}$ , giving

$$\hat{p}_{\text{post}}(0) \approx \frac{1}{\sqrt{2\pi} \hat{s}} \exp\left(-\frac{\hat{\mu}^2}{2\hat{s}^2}\right).$$

- Numerical floors (e.g.,  $10^{-300}$ ) are applied to avoid division by zero or overflow when densities are extremely small.

## 6.5 Interpretation and limitations

- Savage–Dickey requires that the prior under the alternative factorizes and that the null is a point restriction of the alternative.
- The BF depends strongly on the prior scale  $\sigma_{\text{prior}}$ ; different reasonable priors can yield different BFs.
- KDE-based estimates inherit Monte Carlo variability from the posterior samples.

# 7 Bayesian linear regression: model specification and parameterizations

## 7.1 Standard (centered) parameterization

Model:

$$y_i \sim p(y_i | \alpha, \beta, \sigma) = \mathcal{N}(\alpha + \beta x_i, \sigma^2).$$

Priors (example choices used in the code):

$$\alpha \sim \mathcal{N}(0, 5^2), \quad \beta \sim \mathcal{N}(0, \sigma_{\text{prior}}^2), \quad \sigma \sim \text{HalfNormal}(2.5).$$

This is the centered parameterization:  $\beta$  is directly sampled.

## 7.2 Non-centered parameterization

Non-centered parameterization is often used to improve sampling geometry when priors are weakly informative or hierarchical structures exist. Introduce  $\beta_{\text{raw}} \sim \mathcal{N}(0, 1)$  and set

$$\beta = \beta_{\text{raw}} \cdot \sigma_{\text{prior}}.$$

This reparameterization decouples scale from standard normal draws and can reduce funnel-shaped posterior geometries that cause sampler difficulties.

## 7.3 Why non-centered helps

When the posterior is concentrated in a narrow region relative to the prior scale, the centered parameterization can produce strong posterior correlations between  $\beta$  and  $\sigma$  (or other parameters), leading to slow mixing and divergences. The non-centered parameterization often yields more isotropic posterior geometry for the sampler.

# 8 Student– $t$ likelihood for robust regression

## 8.1 Student– $t$ density

The Student– $t$  distribution with location  $\mu$ , scale  $\sigma$ , and degrees of freedom  $\nu > 0$  has density

$$p(y | \mu, \sigma, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})\sqrt{\nu\pi}\sigma} \left(1 + \frac{(y - \mu)^2}{\nu\sigma^2}\right)^{-(\nu+1)/2}.$$

For regression, set  $\mu_i = \alpha + \beta x_i$ .

## 8.2 Robustness properties

- For small  $\nu$  (e.g.,  $\nu \in (1, 10)$ ), the Student– $t$  has heavier tails than the Gaussian, reducing the influence of outliers on parameter estimates.
- As  $\nu \rightarrow \infty$ , the Student– $t$  converges to the Gaussian with variance  $\sigma^2$ .
- The scale parameter  $\sigma$  in the Student– $t$  is not the same as the Gaussian variance; the variance exists only for  $\nu > 2$  and equals  $\sigma^2 \cdot \nu / (\nu - 2)$ .

## 8.3 Priors for $\nu$

A common weakly informative prior for  $\nu$  is an exponential or gamma distribution shifted to favor moderate tail heaviness (e.g.,  $\nu \sim \text{Exponential}(\lambda)$  or  $\nu \sim \text{Exponential}(1/30)$  as in the code). This allows the data to inform tail behavior.

# 9 MCMC sampling diagnostics

Reliable inference from MCMC requires diagnostics. We discuss divergences, Gelman–Rubin  $\hat{R}$ , and effective sample size.

## 9.1 Divergences in Hamiltonian Monte Carlo (NUTS)

- Hamiltonian Monte Carlo (HMC) simulates Hamiltonian dynamics to propose new states. Numerical integrators (e.g., leapfrog) approximate continuous trajectories.
- A *divergence* occurs when the numerical integrator fails to accurately follow the true Hamiltonian trajectory, often due to regions of high curvature in the posterior (e.g., funnels).
- Divergences indicate that the sampler may not be exploring the posterior correctly; zero divergences is the desired outcome.
- Remedies include reparameterization (non-centered), increasing target\_accept (reducing step size), using more tuning, or changing priors/likelihood (e.g., Student– $t$ ).

## 9.2 Gelman–Rubin $\hat{R}$ statistic

**Setup.** Run  $m$  parallel chains, each of length  $n$  after warmup. Let  $\theta_{ij}$  denote the  $j$ -th draw from chain  $i$ . Define per-chain means  $\bar{\theta}_{i..}$  and overall mean  $\bar{\theta}_{...}$

$$\bar{\theta}_{i..} = \frac{1}{n} \sum_{j=1}^n \theta_{ij}, \quad \bar{\theta}_{...} = \frac{1}{m} \sum_{i=1}^m \bar{\theta}_{i..}$$

**Between-chain variance  $B$  and within-chain variance  $W$ .**

$$B = \frac{n}{m-1} \sum_{i=1}^m (\bar{\theta}_{i\cdot} - \bar{\theta}_{..})^2, \quad W = \frac{1}{m} \sum_{i=1}^m s_i^2,$$

where  $s_i^2 = \frac{1}{n-1} \sum_{j=1}^n (\theta_{ij} - \bar{\theta}_{i\cdot})^2$ .

**Estimate of marginal posterior variance**

$$\widehat{\text{Var}}^+ = \frac{n-1}{n} W + \frac{1}{n} B.$$

**Potential scale reduction factor**

$$\hat{R} = \sqrt{\frac{\widehat{\text{Var}}^+}{W}}.$$

Interpretation:  $\hat{R} \rightarrow 1$  as chains converge. Practical thresholds:  $\hat{R} \leq 1.01$  is excellent;  $\hat{R} \leq 1.05$  is often considered acceptable in applied work; values  $> 1.1$  indicate problems.

### 9.3 Effective sample size (ESS)

MCMC draws are autocorrelated; ESS estimates the number of independent samples equivalent to the correlated chain. For a single chain, ESS can be estimated using the autocorrelation function  $\rho_k$ :

$$\text{ESS} \approx \frac{n}{1 + 2 \sum_{k=1}^{\infty} \rho_k}.$$

Arviz and other packages compute multichain ESS using spectral estimators. Larger ESS implies more precise Monte Carlo estimates.

## 10 Posterior predictive checks (PPC)

### 10.1 Posterior predictive distribution

Given posterior draws  $\{\theta^{(s)}\}_{s=1}^S$ , the posterior predictive distribution for a new or replicated observation  $y^{\text{rep}}$  is

$$p(y^{\text{rep}} | y) = \int p(y^{\text{rep}} | \theta) p(\theta | y) d\theta \approx \frac{1}{S} \sum_{s=1}^S p(y^{\text{rep}} | \theta^{(s)}).$$

### 10.2 Generating replicated datasets

For each posterior draw  $\theta^{(s)}$ , simulate  $y^{\text{rep},(s)}$  from the likelihood  $p(y | \theta^{(s)})$ . This yields replicated datasets  $\{y^{\text{rep},(s)}\}$ .

### 10.3 Summary statistics and discrepancy measures

Compute summary statistics  $T(y)$  (e.g., mean, standard deviation) for observed data and for each replicate  $T(y^{\text{rep},(s)})$ . The PPC p-value (one-sided) is estimated as

$$\Pr(T(y^{\text{rep}}) \leq T(y) | y) \approx \frac{1}{S} \sum_{s=1}^S \mathbf{1}\{T(y^{\text{rep},(s)}) \leq T(y)\}.$$

Interpretation: extreme values (close to 0 or 1) indicate systematic misfit for that statistic.

## 11 Leave-one-out cross-validation (LOO) and PSIS

### 11.1 LOO predictive density

The leave-one-out predictive density for observation  $i$  is

$$p(y_i | y_{-i}) = \int p(y_i | \theta) p(\theta | y_{-i}) d\theta.$$

Exact LOO requires refitting the model  $n$  times; approximate methods reuse full-data posterior draws.

### 11.2 Importance sampling identity

Using full-data posterior draws  $\{\theta^{(s)}\}$ , we can write

$$p(y_i | y_{-i}) = \int p(y_i | \theta) \frac{p(\theta | y_{-i})}{p(\theta | y)} p(\theta | y) d\theta = \int p(y_i | \theta) \frac{1}{p(y_i | \theta)} \frac{p(\theta)p(y_{-i} | \theta)}{p(\theta)p(y | \theta)} p(\theta | y) d\theta,$$

which simplifies to an importance sampling estimator with weights proportional to  $1/p(y_i | \theta^{(s)})$ .

### 11.3 Pareto-smoothed importance sampling (PSIS)

Raw importance weights can be unstable. PSIS fits a generalized Pareto distribution to the upper tail of the importance weights and smooths them, producing stabilized weights. The Pareto shape parameter  $k$  diagnoses reliability:

- $k < 0.5$ : reliable.
- $0.5 \leq k < 0.7$ : usable with caution.
- $k \geq 0.7$ : unreliable; LOO estimates may be biased.

Arviz and other libraries compute PSIS-LOO and report Pareto  $k$  values.

### 11.4 LOO model comparison

Compute expected log predictive density (ELPD) or LOO information criterion for each model; prefer the model with higher ELPD (or lower LOOIC). Differences can be assessed with standard errors.

## 12 Highest density interval (HDI) and credible intervals

### 12.1 Definition

A  $100(1 - \alpha)\%$  highest density interval (HDI) for a posterior distribution is the narrowest interval  $[a, b]$  such that

$$\int_a^b p(\beta | y) d\beta = 1 - \alpha$$

and for all  $x \in [a, b]$  and  $x' \notin [a, b]$ ,  $p(x | y) \geq p(x' | y)$ . For unimodal posteriors this is the shortest credible interval.

### 12.2 HDI exclusion of zero

If the 95% HDI does not contain zero (i.e., either  $a > 0$  or  $b < 0$ ), this is often interpreted as evidence that  $\beta$  is credibly nonzero at the 95% level. The code counts how many prior settings yield HDIs excluding zero.

## 13 Combining evidence: heuristics used in the pipeline

The code aggregates multiple evidence sources:

- **Permutation p-value** for Pearson correlation (frequentist).
- **Savage–Dickey BF** estimates across a grid of prior scales; geometric mean of BF values is used to summarize across priors:

$$BF_{geo} = \exp\left(\frac{1}{G} \sum_{g=1}^G \log BF_g\right),$$

where  $G$  is the number of prior settings.

- **BIC approximation BF** as a large-sample check.
- **HDI exclusion counts** across priors.
- **PPC summaries** (proportion of replicates with mean or sd less than or equal to observed).
- **LOO preference for robust model** counts.

A conservative decision rule in the code:

- If  $BF_{geo} \geq 10$ : conclude “no relation” (evidence for null).
- If  $BF_{geo} \leq 0.1$ : conclude “relation” (evidence for alternative).
- If HDI exclusion proportion  $\geq 0.5$  and permutation  $p < 0.05$ : conclude “relation”.
- Otherwise: “inconclusive”.

These thresholds are heuristic and chosen for conservatism.

## 14 Numerical and practical considerations

### 14.1 Finite precision and floors

Densities and likelihoods can be extremely small; numerical floors (e.g.,  $10^{-300}$ ) are used to avoid division by zero and floating point underflow.

### 14.2 Monte Carlo variability

All sample-based estimates (KDE, posterior summaries, PPC proportions, PSIS-LOO) have Monte Carlo error. Increasing draws and chains improves precision.

### 14.3 Standardization of predictors

Standardizing  $x$  (centering and scaling) improves numerical stability and interpretability of priors on  $\beta$ . If  $x$  is standardized to mean 0 and sd 1, a prior  $\beta \sim \mathcal{N}(0, \sigma_{\text{prior}}^2)$  has a direct interpretation in terms of expected change in  $y$  per standard deviation change in  $x$ .

## 15 Worked symbolic derivations and small proofs

### 15.1 Derivation: relation between $r$ and OLS slope

Starting from OLS slope:

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}}.$$

Divide numerator and denominator by  $\sqrt{S_{xx}S_{yy}}$ :

$$\hat{\beta} = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} \cdot \frac{\sqrt{S_{xx}S_{yy}}}{S_{xx}} = r \cdot \frac{\sqrt{S_{yy}}}{\sqrt{S_{xx}}}.$$

Thus  $r = \hat{\beta} \cdot \sqrt{S_{xx}/S_{yy}}$ .

### 15.2 Derivation: BIC for Gaussian linear model

For Gaussian errors with variance  $\sigma^2$ , the log-likelihood is

$$\log L(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{RSS}.$$

Maximizing over  $\sigma^2$  yields  $\hat{\sigma}^2 = \text{RSS}/n$ . Plugging in,

$$\log L(\hat{\theta}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log\left(\frac{\text{RSS}}{n}\right) - \frac{n}{2}.$$

Up to constants independent of model,  $-2 \log L(\hat{\theta})$  is proportional to  $n \log(\text{RSS}/n)$ . Adding the penalty  $k \log n$  yields the BIC formula used earlier.

### 15.3 Derivation: Savage–Dickey identity (sketch)

Let  $M_1$  be the full model with parameter  $\beta$  and nuisance  $\phi$ . The marginal posterior of  $\beta$  under  $M_1$  is

$$p(\beta \mid y, M_1) = \frac{p(\beta) \int p(y \mid \beta, \phi)p(\phi) d\phi}{p(y \mid M_1)}.$$

Evaluating at  $\beta = 0$ ,

$$p(0 \mid y, M_1) = \frac{p(0) \int p(y \mid 0, \phi)p(\phi) d\phi}{p(y \mid M_1)}.$$

But the numerator integral is exactly  $p(y \mid M_0)$  (the marginal likelihood under the null). Rearranging gives

$$\frac{p(y \mid M_0)}{p(y \mid M_1)} = \frac{p(0)}{p(0 \mid y, M_1)}.$$

This is the Savage–Dickey identity.

## 16 Summary of recommended practical workflow (algorithmic steps)

1. **Data cleaning and alignment.** Remove non-finite values and align paired observations.
2. **Frequentist baseline.** Fit OLS, compute slope,  $R^2$ , Pearson  $r$ , and permutation p-value with  $N_{\text{perm}}$  permutations.
3. **Bayesian grid.** For a grid of prior scales  $\{\sigma_{\text{prior},g}\}$ , fit Bayesian linear models using non-centered parameterization and optionally Student– $t$  likelihoods if sampling issues arise.
4. **Diagnostics.** For each fit, check divergences (should be 0),  $\hat{R}$  (should be near 1), and ESS.
5. **Posterior summaries.** Extract posterior samples of  $\beta$ , compute means, sds, and 95% HDIs.
6. **Savage–Dickey BF.** Estimate posterior density at 0 via KDE (or normal approximation fallback) and compute BF for each prior scale.
7. **Aggregate BF.** Compute geometric mean of BF values across priors to summarize sensitivity to prior scale.
8. **PPC and LOO.** Compute posterior predictive summaries and PSIS-LOO comparisons between Gaussian and robust (Student– $t$ ) fits.
9. **Decision heuristic.** Combine BF, HDI counts, permutation p-value, PPC, and LOO preferences to form a conservative verdict.

## 17 Appendix: useful formulae and constants

$$\text{Normal density: } \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right).$$

$$\text{Student-}t \text{ density: } t(x | \mu, \sigma, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})\sqrt{\nu\pi}\sigma} \left(1 + \frac{(x - \mu)^2}{\nu\sigma^2}\right)^{-(\nu+1)/2}.$$

$$\text{Silverman bandwidth (rule of thumb): } h = 0.9 \min\{\hat{\sigma}, \text{IQR}/1.34\} S^{-1/5}.$$

## Concluding remarks

This document has integrated and expanded every mathematical element present in the analysis pipeline, providing definitions, derivations, practical estimation formulas, and numerical considerations. The pipeline intentionally combines frequentist permutation inference, asymptotic BIC approximations, and sample-based Bayesian evidence (Savage–Dickey) together with robust modeling and predictive checks to produce a conservative, multi-angle assessment of evidence for or against a linear relationship between two series.