# A Bayesian approach to asteroid hazard classification: Logit vs Probit

Alex Costa, Sara Pascali | Università Cattolica del Sacro Cuore | Bayesian Modelling

25/06/2024

# 1 Introduction

In this project we adopt a **Bayesian regression** framework to model the probability that an asteroid poses an hazard to Earth, incorporating prior beliefs and accounting for parameter uncertainty. Understanding the key orbital and physical characteristics that contribute to asteroid hazard classification is critical for planetary defense strategies and space mission planning. The dataset used is publicly available on Kaggle (sourced from NASA's Center for Near Earth Object Studies): https://www.kaggle.com/datasets/lovishbansal123/nasa-asteroids-classification

The Bayesian approach is particularly suited for this context due to the interpretability of posterior distributions, the ability to incorporate prior expert knowledge, and the robustness it offers in small-sample or imbalanced data scenarios. In our dataset indeed the response variable **Hazardous** is notably unbalanced, with only about 15.5% of the 4,687 asteroids classified as hazardous.

We begin with exploratory data analysis and preprocessing, followed by the implementation of Bayesian logistic and Bayesian probit regression models. Model performance is evaluated through posterior predictive checks, classification accuracy, and ROC curve analysis on a held-out test set.

In addition to the response variable, the dataset contains 15 predictors:

- **Absolute Magnitude**: A standardized measure of an asteroid's luminosity.

- **Estimated Diameter**: An estimate of the diameter of an asteroid expressed in kilometers.

- **Relative Velocity**: The speed at which an asteroid moves relative to Earth during its close approach, expressed in kilometers per second.

- **Miss Distance**: The minimum distance at which the asteroid will pass by Earth during a close approach, expressed in kilometers.

- **Orbit Uncertainty**: The uncertainty of the estimated orbital parameters: 1 = Low Uncertainty; 2 = Medium Uncertainty; 3 = High Uncertainty.

- **Minimum Orbit Intersection**: The minimum distance between the object's orbit and Earth's orbit, expressed in Astronomical Units (1 AU = 1,496e+8 Km).

- **Eccentricity**: The shape of the orbit, with values ranging from 0 (circular orbit) to 1 (highly elongated elliptical orbit).

- **Semi Major Axis**: One half of the longest diameter of the elliptical orbit, expressed in Astronomical Units. It essentially defines the size of the orbit.

- **Inclination**: The angle between the asteroid's orbital plane and the ecliptic plane (Earth's orbital plane), evaluated in degrees.

- **Orbital Period**: The time required for the asteroid to complete one full orbit around the Sun, expressed in years.

- **Perihelion Distance**: The distance from the asteroid's Perihelion (the closest point in the orbit to the Sun) to the Sun, expressed in Astronomical Units.

- **Perihelion Argument**: The angle between the ascending node and the Perihelion, expressed in degrees.

- **Aphelion Distance**: The distance from the asteroid's Aphelion (the farthest point in the orbit to the Sun) to the Sun, expressed in Astronomical Units.

- **Mean Anomaly**: The fraction of the orbit already completed by the asteroid since Perihelion, expressed in degrees.

- **Mean Motion**: The average angular speed at which the asteroid orbits the Sun, expressed in degrees per day.

## 2. Exploratory Data Analysis

To evaluate model performance and ensure generalizability, we partition the dataset of 4,687 asteroid observations into two subsets: 80% of the data is allocated to the training set, which is used for model fitting, while the remaining 20% forms the test set, reserved for out-of-sample evaluation. This split allows us to assess how well the Bayesian models perform on unseen data.

We start with an exploratory analysis of the training data, presenting the most relevant plots below:
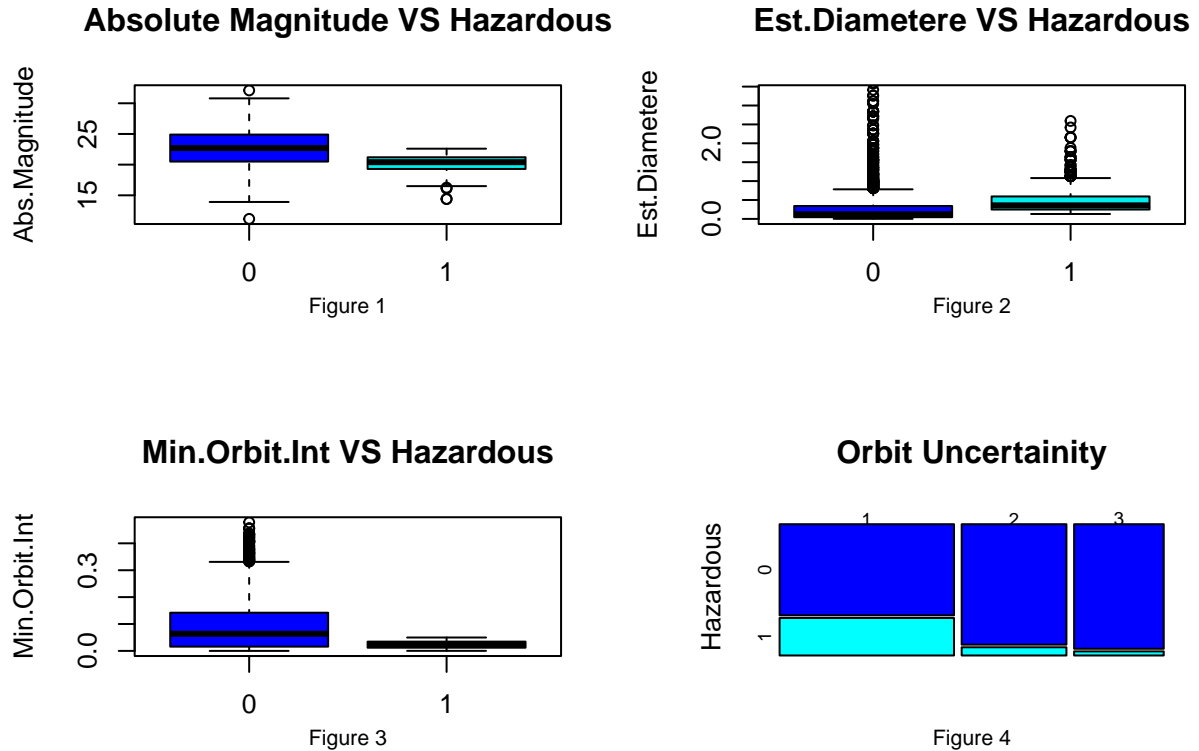


Figure 1: hazardous asteroids tend to have lower absolute magnitudes, meaning they are brighter or more luminous. Since brightness often correlates with size and visibility, this might suggest hazardous asteroids are larger or more detectable.

FIgure 2: hazardous asteroids generally have larger estimated diameters, but the distribution is highly skewed (with some extreme outliers).

Figure 3: hazardous asteroids have a significantly lower Minimum Orbit Intersection, meaning they pass closer to Earth's orbit.

Figure 4: almost all hazardous asteroids have an orbit uncertainty level 1, suggesting that hazardous asteroids tend to have well-defined orbits, likely because they are tracked more carefully.

Prior to fitting the Bayesian models, we standardize the numeric predictors in both the training and test sets by subtracting the sample mean and dividing by the standard deviation (z-score standardization). This is applied to all continuous variables, excluding the intercept and categorical predictors. Standardization enhances the efficiency and convergence of Markov Chain Monte Carlo (MCMC) algorithms and enables a more interpretable comparison of the resulting regression coefficients, as they are placed on a common scale.

# 3 GLM for binary response variables

Generalized Linear Models (GLMs) extend ordinary linear regression to response variables that follow distributions other than the Normal.

In our context, the response variable is binary: it takes value 1 if the asteroid is classified as hazardous, and 0 otherwise. Therefore, the appropriate likelihood is a Bernoulli distribution with success probability $\pi_i$ representing the probability that the $i$-th asteroid is hazardous.

GLMs link the expected value of the response to a linear combination of the predictors via a **link function**. In this analysis, we consider and compare two popular link functions for binary responses: the logit and the probit.

## 3.1 Logistic regression

Logistic regression uses the logit link function:

$$\pi_i = h(\boldsymbol{\beta}^T \boldsymbol{x}_i) = \frac{e^{\boldsymbol{\beta}^T \boldsymbol{x}_i}}{1 + e^{\boldsymbol{\beta}^T \boldsymbol{x}_i}} \qquad \Leftrightarrow \qquad log(\frac{\pi_i}{1 - \pi_i}) = \boldsymbol{\beta}^T \boldsymbol{x}_i$$

The corresponding likelihood function is:

$$p(\boldsymbol{y}|\boldsymbol{\beta}) = \prod_{i=1}^{n} p(y_i|\pi_i) = \prod_{i=1}^{n} h(\boldsymbol{\beta}^T \boldsymbol{x}_i)^{y_i} (1 - h(\boldsymbol{\beta}^T \boldsymbol{x}_i))^{1-y_i}$$

We place **independent Normal** priors on the regression coefficients, so that the joint prior is the product of the marginal priors:

$$\beta_j \overset{ind}{\sim} N(\beta_{0j}, \sigma_{0j}^2) \quad \Rightarrow \quad p(\boldsymbol{\beta}) = \prod_{j=1}^{p} dN(\beta_j|\beta_{0j}, \sigma_{0j}^2)$$

Since this prior is not conjugate, we estimate the posterior distribution via the **Metropolis-Hastings algorithm**.

## 3.2 Probit regression

In probit regression, the link function is the inverse of the standard normal CDF:

$$\pi_i = h(\boldsymbol{\beta}^T \boldsymbol{x}_i) = \Phi(\boldsymbol{\beta}^T \boldsymbol{x}_i) \qquad \Leftrightarrow \qquad \Phi^{-1}(\pi_i) = \boldsymbol{\beta}^T \boldsymbol{x}_i$$

The likelihood and priors remain the same as in logistic regression. However, the probit model admits a latent variable formulation, allowing us to derive full conditional distributions and implement a **Gibbs sampler**.

Assume that for each observation $y_i$, there exists a latent variable $z_i$ such that:

$$z_i | \boldsymbol{\beta} \overset{ind}{\sim} N(\boldsymbol{\beta}^T \boldsymbol{x}_i, 1) \qquad\qquad \boldsymbol{\beta} \sim N_p(\boldsymbol{\beta}_0, \boldsymbol{V}^{-1})$$

We obtain the following full conditional distributions for $\boldsymbol{z}$ and $\boldsymbol{\beta}$:

$$z_i | \boldsymbol{\beta}, y_i \overset{ind}{\sim} tN(\boldsymbol{\beta}^T \boldsymbol{x}_i, 1, \theta_{y_i-1}, \theta_{y_i}) \qquad\qquad \boldsymbol{\beta} | \boldsymbol{y}, \boldsymbol{z} \sim N_p(\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{V}}^{-1})$$

Sampling alternately from these distributions yields draws from the posterior of $\boldsymbol{\beta}$.

## 3.3 Prediction via Bayesian Model Averaging

Rather than selecting a single model, we use **Bayesian Model Averaging (BMA)** to account for model uncertainty in the predictions. Let $\boldsymbol{\gamma} = (\gamma_1, ..., \gamma_p)^T$ be a binary vector indicating inclusion of each predictor:

$$\gamma_j = \begin{cases} 1 & \text{if } X_j \text{ is included in the model} \\ 0 & \text{otherwise} \end{cases}$$

The model becomes:

$$E(Y|x) = h(\gamma_1 \beta_1 X_1 + ... + \gamma_p \beta_p X_p)$$

We assign the following priors:

$$\beta_j | \gamma_j = 1 \overset{ind}{\sim} N(\beta_{0j}, \sigma_{0j}^2) \qquad\qquad \gamma_j \overset{iid}{\sim} Ber(w) \qquad\qquad w \sim Beta(a, b)$$

This leads to the **spike-and-slab** prior:

$$p(\beta_j, \gamma_j) = (1 - w)\delta_0 + w\, dN(\beta_j | \beta_{0j}, \sigma_{0j}^2)$$

Given $S$ posterior samples $(\boldsymbol{\beta}^{(s)}, \boldsymbol{\gamma}^{(s)})_{s=1}^S$ obtained through Gibbs Sampling or Metropolis Hastings, the predictive distribution for a new input $\boldsymbol{x}^*$ is approximated by:

1) Compute $\eta^{(s)} = \gamma_1^{(s)} \beta_1^{(s)} x_1^* + ... + \gamma_p^{(s)} \beta_p^{(s)} x_p^*$

2) Compute $\mu^{(s)} = h(\eta^{(s)})$

3) Sample $y^{*(s)}$ from $p(y^* | \mu^{(s)})$

The result is a posterior predictive sample for the new observation, capturing uncertainty in both parameter values and model structure.

# 4 Posterior distributions approximation

## 4.1 Logit regression

We apply the logistic regression model as introduced in Section 3.1 to the asteroid dataset. Since the model does not admit a closed-form posterior distribution, we use the Metropolis-Hastings algorithm to sample from the posterior of the regression coefficients $\beta_j$.

All numeric predictors are standardized, allowing us to adopt weakly informative priors centered at zero with unit variance:

$$\beta_j \sim N(0, 1) \qquad \forall j$$

To perform Bayesian variable selection, we also introduce latent binary indicators $\gamma_j \sim \text{Bern}(w)$, where $w \sim \text{Beta}(1, 1)$ is an unknown prior inclusion probability. The posterior is approximated using a JAGS model. We initialize the sampler with beta = 0 and all variables included (gamma = 1), and run the model for 5,000 iterations with a burn-in of 1,000. Posterior samples are extracted for both $\beta$ and $\gamma$.

**Posterior distribution of Est. Diameter**



Est. Diam.

Figure 5

**Posterior distribution of Min. Orbit Intersectic**



Min. Orbit Int.
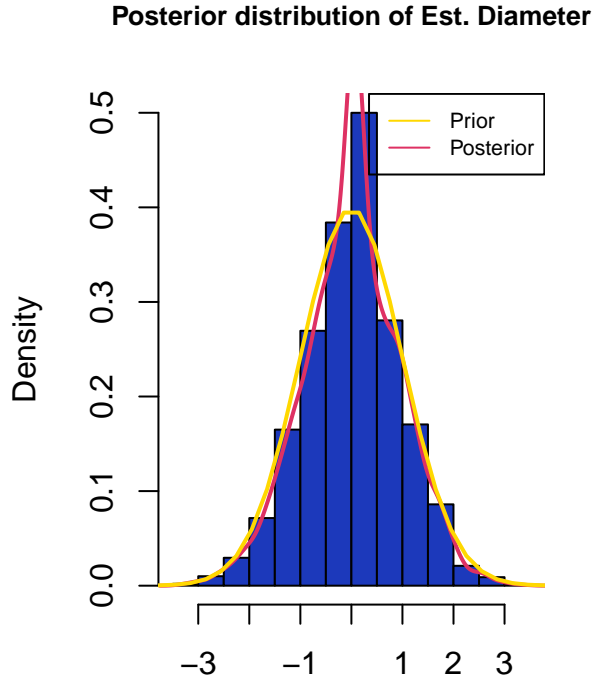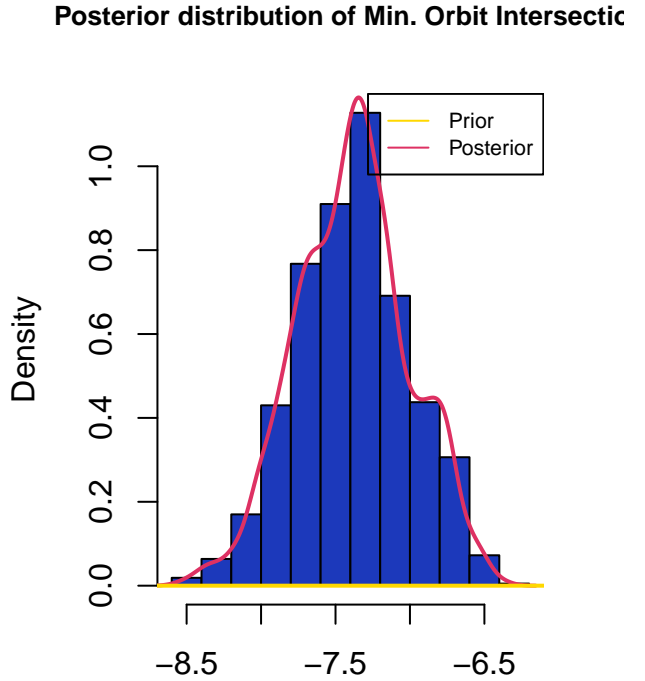
Figure 6

Figures 5 and 6 show the posterior distributions of the coefficients associated with Estimated Diameter and Minimum Orbit Intersection under the logit model, while the red curves are the estimated posterior densities and the yellow lines correspond to the standard normal priors. Compared to the priors, both posterior distributions are more concentrated, indicating that the data have significantly informed the estimation of these coefficients. The coefficient for Estimated Diameter remains centered around zero but with reduced variance, suggesting a moderate and possibly weak effect. In contrast, the posterior distribution for Minimum Orbit Intersection is clearly shifted to the left and displays a sharper peak, reflecting a stronger and more consistent influence on the probability of an asteroid being hazardous.

We compute **summary statistics** for each regression coefficient $\beta_j$, including:

- Posterior Mean

- Posterior Standard Deviation

- 95% Credible Interval

- Posterior Inclusion Probability

Variables like Estimated Diameter and Minimum Orbit Intersection show high inclusion probabilities, meaning they are consistently selected across posterior draws. Others like Miss Distance and Inclination have low inclusion probabilities, suggesting weak or no association with the hazard classification outcome.

To assess the quality of the approximation obtained through the Metropolis-Hastings algorithm, we perform a set of diagnostic checks using the following tools:

**Trace plot**: graphical representation of the sampled values $\theta^{(s)}$ across iteration $s = 1, \ldots, S$. For a good approximation, the chain should exhibit no trends and it should be concentrated within a region of high posterior probability, centered around the mode of $P(\theta|\boldsymbol{y})$.
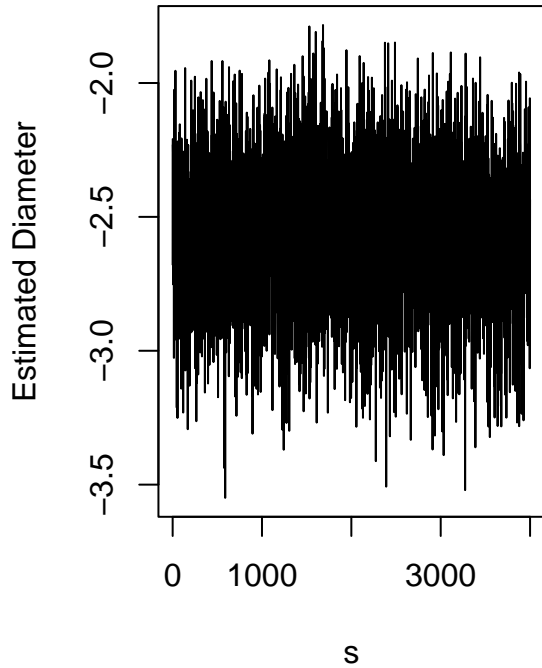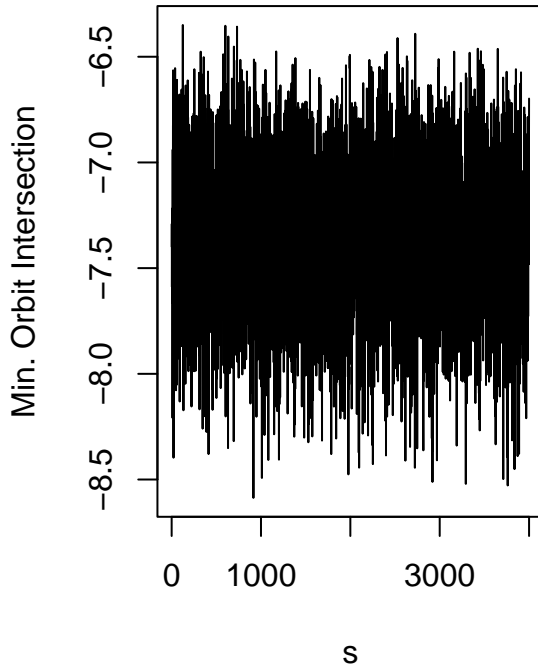


Figure 7



Figure 8

Figures 7 and 8 show the trace plots of the draws for the coefficients associated with the variables Estimated Diameter and Minimum Orbit Intersection. Both plots look nice, indeed there seem to be no trends and the values are quite concentrated in a region, even if there are some peaks. Furthermore, the burn-in period of 1000 draws is sufficient to reach convergence of the chain.

**Auto Correlation Function**: a measure of correlation in a Markov Chain. For a good approximation, the autocorrelation between various lags should decay rapidly towards zero. This indicates that successive observations are weakly correlated.

**Estimated Diameter**
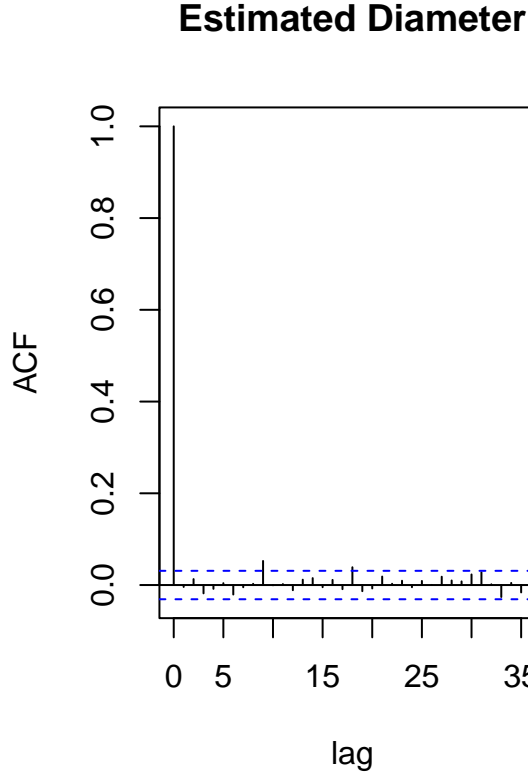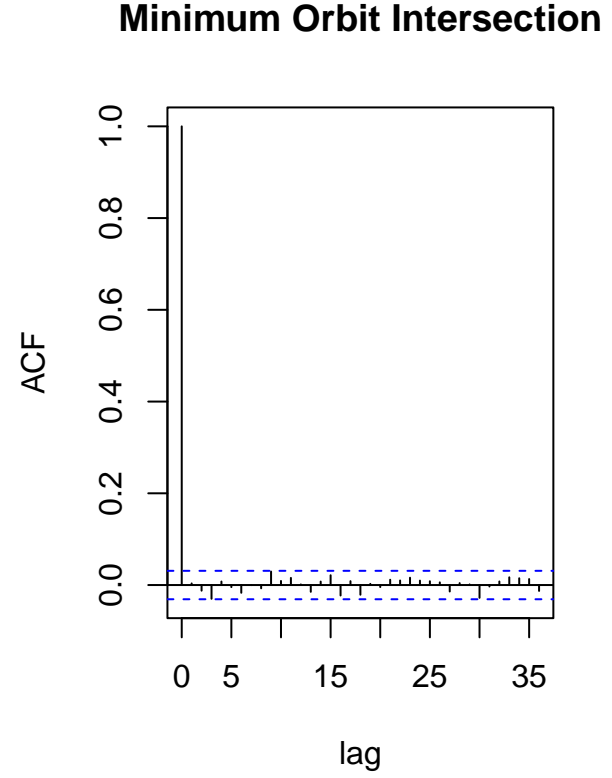
**Minimum Orbit Intersection**



Figure 9



Figure 10

We can see from the ACF for Estimated Diameter (Figure 9) and Minimum Orbit Intersection (Figure 10) that the level of autocorrelation is very low for every lag, and this indicates that we don't need to thin the chain. It holds also for all the other coefficients.

**Geweke test**: a statistical test that compares the means of two subsets of the chain, one consisting of the first X% samples and the other consisting of the last Y% samples. If the chain converges, it is expected that these means are similar. Otherwise we should consider increasing the burn-in period. We choose X = Y = 10%.

| beta[1] | beta[2] | beta[3] | beta[4] | beta[5] | beta[6] | beta[7] | beta[8] | beta[9] |
|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| 0.354 | 1.4158 | 1.2709 | -1.2848 | 1.1373 | -0.1867 | 1.4445 | 0.5275 | -0.1276 |

| beta[10] | beta[11] | beta[12] | beta[13] | beta[14] | beta[15] | beta[16] | beta[17] |
|----------|----------|----------|----------|----------|----------|----------|----------|
| 0.5072 | -0.3676 | -0.0578 | 0.2064 | -1.587 | -0.9779 | 0.3258 | 0.5478 |

The output of the test are the values of the test statistic for each of the parameters in the model. Such statistic is asymptotically Standard Normal, therefore at significance level 5% we don't reject the null hypothesis for any $\beta$, since all the values are between -1.96 and 1.96. Thus we can say that all the chains have reached convergence and the burn-in period of 1000 draws is adequate.

**Effective Sample Size**: measures how much information is loss due to autocorrelation in the sequence. Although our sequence may have a sample size of N, our effective sample size could be smaller or even bigger due to the correlation and redundancy between the samples.
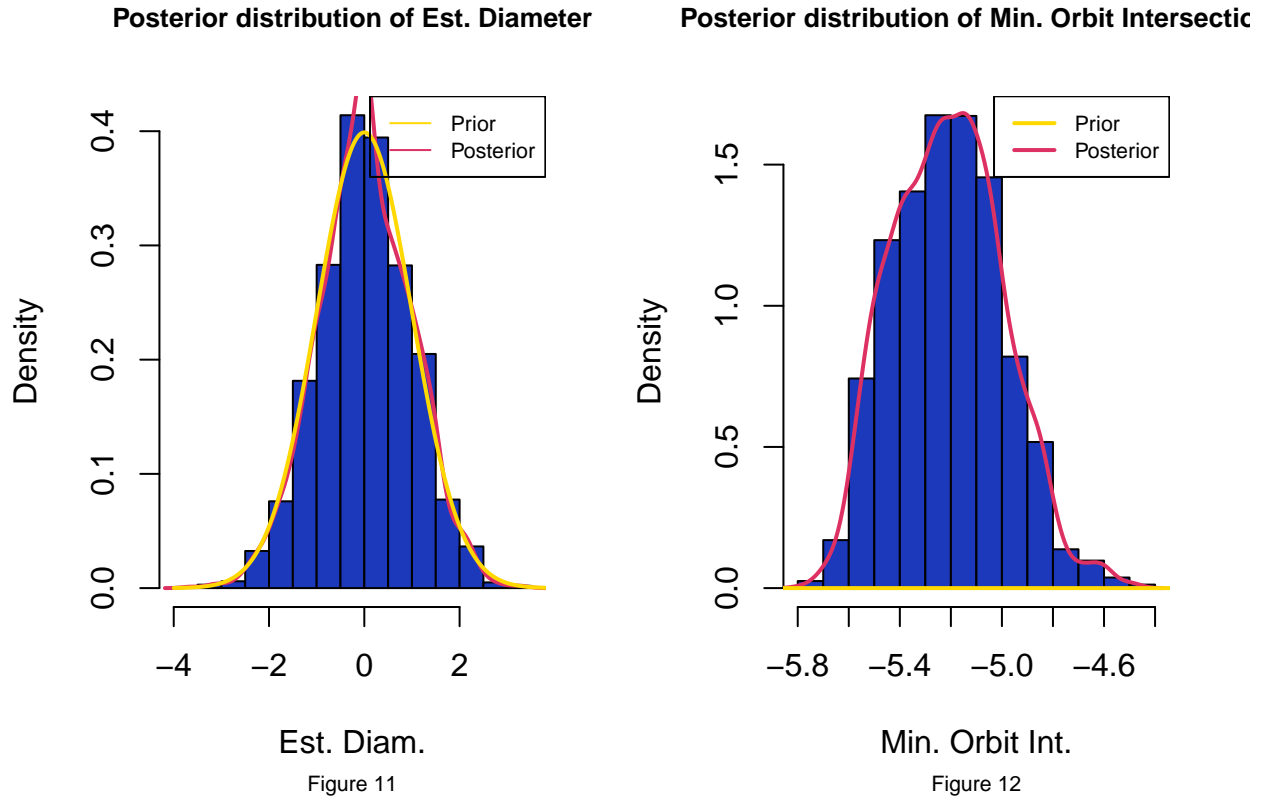
| beta[1] | beta[2] | beta[3] | beta[4] | beta[5] | beta[6] | beta[7] | beta[8] | beta[9] |
|---|---|---|---|---|---|---|---|---|
| 3819.541 | 4000 | 4000 | 4000 | 4000 | 4000 | 4000 | 4000 | 4000 |

| beta[10] | beta[11] | beta[12] | beta[13] | beta[14] | beta[15] | beta[16] | beta[17] |
|---|---|---|---|---|---|---|---|
| 4000 | 4000 | 4000 | 4000 | 4000 | 4000 | 4000 | 4000 |

All the estimated coefficients have an effective sample size equal or almost equal to the real sample size, 4000. This confirms that we don't have autocorrelation issues.

## 4.2 Probit regression

We apply the probit regression model, as outlined in Section 2.2, to the asteroid dataset using a Metropolis-Hastings algorithm, without relying on the latent variable representation. The hyperparameters for the priors are set to the same values used in the previous chapter.

**Posterior distribution of Est. Diameter**



Est. Diam.

Figure 11

**Posterior distribution of Min. Orbit Intersectic**



Min. Orbit Int.

Figure 12

Figures 11 and 12 display the posterior distributions of the coefficients for the variables Estimated Diameter and Minimum Orbit Intersection, respectively, obtained under the probit model, while the yellow curves show the standard normal priors, and the red curves represent the posterior densities. In both cases, we observe a very similar behavior with respect to the logit model: the posterior distribution for Estimated Diameter remains centered around zero but with reduced variance, suggesting a moderate and possibly weak effect, while the posterior distribution for Minimum Orbit Intersection is clearly shifted to the left and displays a sharper peak.

We check again the same **diagnostics** tools, reporting the results only for one coefficient, since the results are very similar to the previous paragraph:
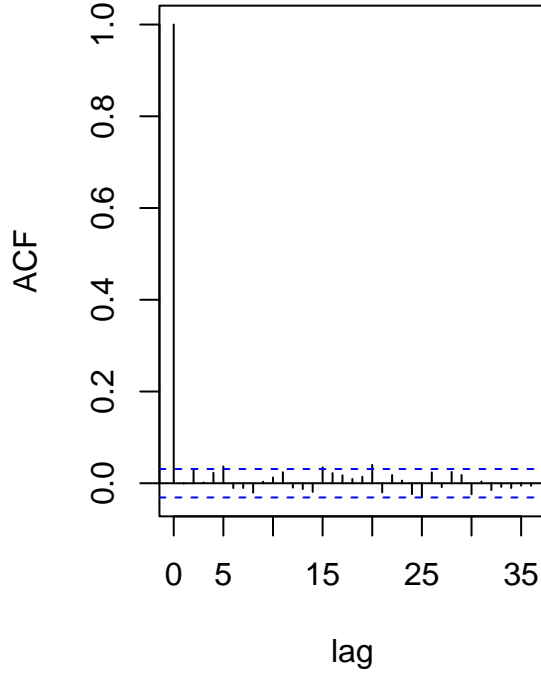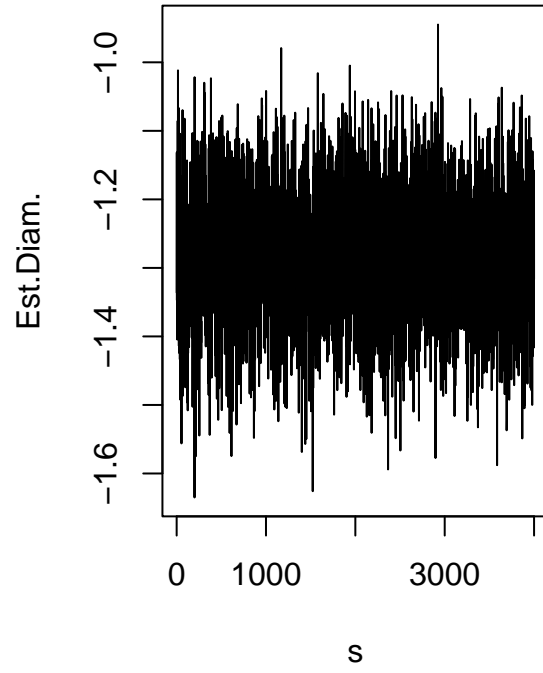
## Estimated Diameter



Figure 13



Figure 14

Figures 13 and 14 display the autocorrelation function (ACF) and the trace plot for the posterior draws of the Estimated Diameter coefficient in the probit model. In Figure 13, the ACF drops rapidly and remains close to zero across lags, indicating low autocorrelation among the sampled values and thus good mixing of the Markov chain. Figure 14 further supports this by showing a stable trace plot with no evident trends or drift, suggesting that the chain has reached convergence and is exploring the posterior distribution efficiently. The absence of significant autocorrelation and the dense concentration of the chain around a fixed region confirm the reliability of the samples for inference.

**Geweke test:**

| beta[1] | beta[2] | beta[3] | beta[4] | beta[5] | beta[6] | beta[7] | beta[8] | beta[9] |
|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| 0.354 | 1.4158 | 1.2709 | -1.2848 | 1.1373 | -0.1867 | 1.4445 | 0.5275 | -0.1276 |

| beta[10] | beta[11] | beta[12] | beta[13] | beta[14] | beta[15] | beta[16] | beta[17] |
|----------|----------|----------|----------|----------|----------|----------|----------|
| 0.5072 | -0.3676 | -0.0578 | 0.2064 | -1.587 | -0.9779 | 0.3258 | 0.5478 |

For the Geweke test, again we don't reject the null hypothesis for any $\beta$ at a 5% significance level.
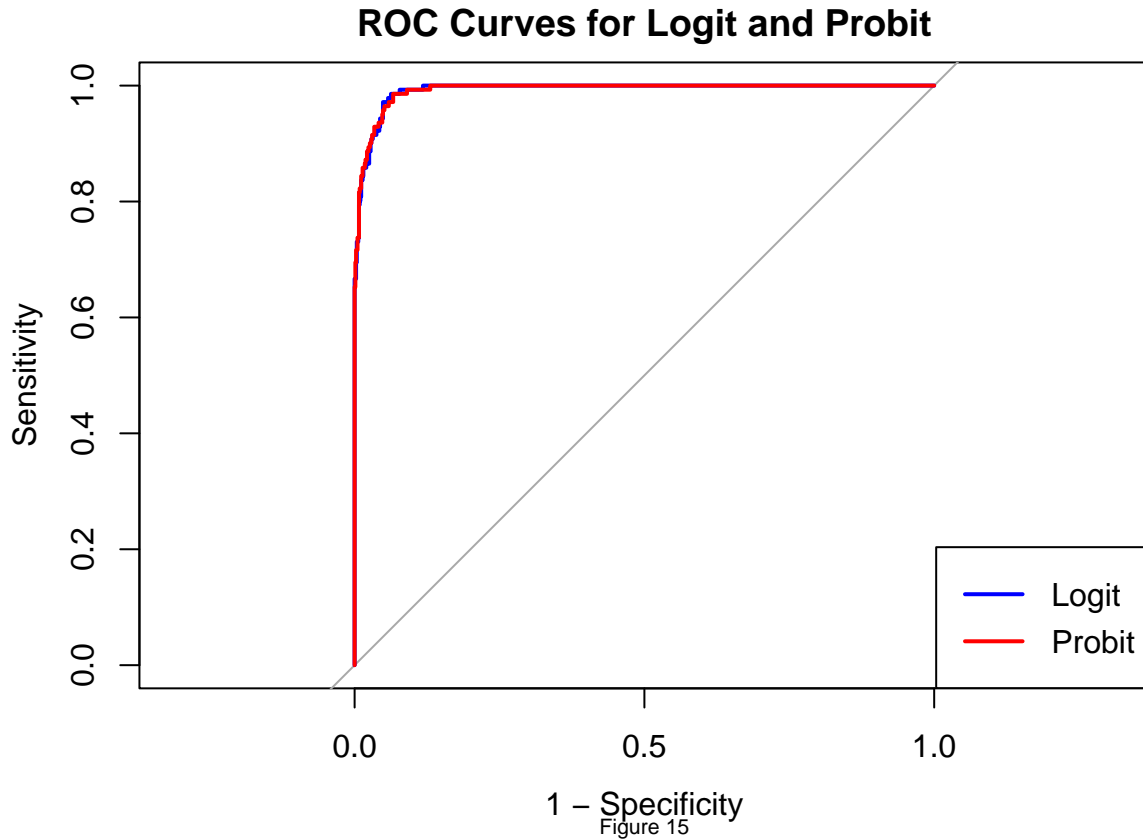
**Effective Sample Size:**

9

| beta[1] | beta[2] | beta[3] | beta[4] | beta[5] | beta[6] | beta[7] | beta[8] | beta[9] |
|---|---|---|---|---|---|---|---|---|
| 4000 | 4000 | 3356.113 | 4000 | 4000 | 4000 | 4000 | 3057.136 | 3604.868 |

| beta[10] | beta[11] | beta[12] | beta[13] | beta[14] | beta[15] | beta[16] | beta[17] |
|---|---|---|---|---|---|---|---|
| 4285.356 | 4000 | 4000 | 4000 | 4389.372 | 3873.549 | 4000 | 4000 |

In general, the ESS values are very similar or even slightly exceeding the nominal sample size of 4000, indicating excellent mixing of the Markov chain and a low level of autocorrelation. The few coefficients with slightly lower ESS values still retain sufficient effective sample sizes to ensure reliable posterior estimates.

# 5 Prediction: comparison between Logit and Probit models

We now evaluate the predictive performance of the logit and probit regression models on the test dataset. Using Bayesian Model Averaging (BMA), we estimate the probability that each asteroid is hazardous. These probabilities are then converted into binary predictions based on a classification threshold. We select the best threshold as the one that maximizes the test accuracy. Then, we assess the best model performance using a **confusion matrix**.



Figure 15

The ROC curve visually confirms the strong performance of both models. The curves for both models almost completely overlap, meaning they are nearly indistinguishable in their ability to separate the two classes

```
## Best threshold (Logit): 0.53
```

10

```
## Test Accuracy (Logit): 0.9658849

## AUC (Logit): 0.9922404

## Best threshold (Probit): 0.54

## Test Accuracy (Probit): 0.9658849

## AUC (Probit): 0.992276
```

Despite slight differences in thresholds and AUC values, both models achieve the same test accuracy, indicating that in practical terms, either model would be equally effective. We show next the confusion matrix and the classification report for the probit model with the best threshold:

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0    1
##          0 785   20
##          1  12  121
##
##                Accuracy : 0.9659
##                  95% CI : (0.9522, 0.9766)
##     No Information Rate : 0.8497
##     P-Value [Acc > NIR] : <2e-16
##
##                   Kappa : 0.8633
##
##  Mcnemar's Test P-Value : 0.2159
##
##             Sensitivity : 0.8582
##             Specificity : 0.9849
##          Pos Pred Value : 0.9098
##          Neg Pred Value : 0.9752
##              Prevalence : 0.1503
##          Detection Rate : 0.1290
##    Detection Prevalence : 0.1418
##       Balanced Accuracy : 0.9215
##
##        'Positive' Class : 1
##
```

The high values of Accuracy, Sensitivity and Specificity confirm that the probit model offers a very good classification performance at the selected threshold.

# 6 Conclusion

In this project, we applied Bayesian Model Averaging (BMA) with logit and probit models to predict whether an asteroid is classified as hazardous based on its physical and orbital features. Through a comprehensive Bayesian framework, we were able to not only incorporate model uncertainty but also derive full posterior distributions over parameters, allowing for deeper inference and robust predictions.

We thoroughly examined the posterior distributions, assessed convergence diagnostics, and validated model performance using the test set. Key findings include:

- Both the logit and probit models achieved exceptionally high **predictive accuracy** on the test data, with an optimal threshold accuracy of 96.6% and AUC values above 0.99, indicating excellent discriminative ability.

- Posterior and prior comparisons showed that the data had a strong influence on the estimates of key variables such as Estimated Diameter and Minimum Orbit Intersection, which were also supported by effective sample sizes and trace diagnostics.

The results suggest that both modeling approaches are highly suitable for this type of binary classification task in a space science context and they are practically **interchangeable** in terms of performance.

Future extensions of this work could include the use of hierarchical priors, nonlinear predictors, or exploring Bayesian neural networks to capture more complex relationships within the data.