# The Use of Bayesian Techniques with Binary and Vector Data

## By

## Nahed talaat  Abd El kareem

### Doctorate of Philosophy in Applied Statistics, Benha University

**Abstract:**

Bayesian methods are a powerful and flexible class of statistical techniques which can be used to solve a large number of problems, In this research, we discuss the use of the Bayesian  method  in classification problems with binary and vector data

The following results were reached in this research, Bayesian methods are a powerful tool for classification problems with binary and vector data. These methods allow us to incorporate prior knowledge, handle missing data, and quantify uncertainty in our predictions. For binary data, Bayesian methods can be used to model the probability of each class given the observed binary values and make predictions for new observations. For vector data, Bayesian methods can be used to model the probability of each class given the observed vector values and make predictions for new observations.


**Key words:** Bayesian methods, binary data, vector data.

## 1. Introduction:

Bayes methods are a group of statistical methods that use Bayes' theory In order to generate a probability distribution for the parameter x, we utilize the existing information we have regarding this parameter. In Bayesian methods, probabilities are interpreted as degrees of belief, and the goal is to infer the most likely hypothesis given the available data.

There are a wide range of problems to which Bayesian methods can be applied[Albert, 2009] including classification problems with binary and vector data. In the case of binary data, each observation is represented by a vector of binary values, where each value corresponds to the presence or absence of a particular feature. Bayesian methods can be used to model the probability of each class given the observed binary values, and to make predictions for new observations.

In the case of vector data [Gelman, &al 2013] each observation is represented by a vector of continuous or discrete values, where each value corresponds to a particular feature. Bayesian methods can be used to model the probability of each class given the observed vector values, and to make predictions for new observations.

Bayesian methods have several advantages over other statistical techniques, including the ability to incorporate prior knowledge, to handle missing data and to quantify uncertainty. However, Bayesian methods can also be computationally intensive and may require the specification of prior distributions, which can be subjective and difficult to choose.

In general, Bayesian methods are a versatile and robust set of statistical techniques that can effectively address a wide range of problems, for example classification problems with binary and vector data.

## 2. Conceptual framework for research

### 2.1 BAYESIAN ANALYSIS:

In many cases [Barber, 2012] we have additional information from our previous experiences about parameter θ . We may notice that it takes different values and that there is evidence that $\theta$ changes and that this change and additional information can be represented through a probability distribution $\pi(\theta)$ for parameter $\theta$. In other words, the parameter θ is treated as a random variable that follows a probability distribution denoted as π(θ).

Thus, a difference between that traditional statistical methods and the Bayesian method [McElreath,2020] that means that the parameter θ is regarded as a random variable that is characterized by π(θ) which represents a probability distribution that represents our prior knowledge or information about θ and describes the degree of our belief in the possible values of this parameter or describes our previous experience about the parameter before obtaining the sample, and accordingly this distribution is called prior distribution: It describes the information and past experience that we have about parameter $\theta$.

## 2.2 The prior distribution:

In bayesian analysis, the prior distribution [McElreath,2020]  π(θ) for parameter θ represents our initial knowledge or beliefs about the parameters of interest before we have any data. It allows us to quantify our uncertainty about the parameters before we take the data into account.

The prior distribution is specified based on prior knowledge, previous studies, expert opinions, or subjective beliefs. It can take various forms, such as a normal distribution, uniform distribution, beta distribution, or any other probability distribution that is appropriate for the parameter being modeled.

The selection of the prior distribution [Gelman, &al 2013]. can greatly influence the posterior distribution and the subsequent conclusions drawn from it. A prior can be informative, where it assigns relatively higher probability to certain values of the parameter based on strong prior knowledge, or it can be non-informative, where it assigns relatively equal probability to a wide range of values. Non-informative priors are often used when there is limited prior knowledge or when we want the data to dominate the inference.

By applying Bayes' theorem, the prior distribution is combined with the likelihood function, which indicates the probability of observing the data given the parameters. This combination results in the posterior distribution, which reflects the revised knowledge about the parameters after taking the data into consideration.

One of the pros of the Bayesian approach is that it allows for the iterative update of the prior distribution as more data becomes available. This is particularly useful when dealing with sequential or streaming data.

It is worth noting that the selection of the prior distribution can be subjective, and different individuals might hold different prior beliefs. To evaluate the influence of different prior specifications on the outcomes, sensitivity analysis or robustness checks can be conducted.

## 2.3 Posterior distribution:

The posterior distribution, as described by [McElreath,2020], is the revised probability distribution of the parameters of interest. It is obtained by combining the prior distribution with the likelihood function using Bayes' theorem, after incorporating the information contained in the observed data.

Mathematically, the posterior distribution[Robert, 2007] is calculated as:

Posterior distribution $\propto$ Prior distribution $\times$ Likelihood function

$$\pi(\theta/x) \propto \pi(\theta)l(\theta/x). \quad (1)$$

The posterior distribution reflects our revised understanding of the parameters based on the observed data. It provides a complete probability distribution that reflects the uncertainty in the parameter estimates.

## 2.4 Binary data:

Binary data [Hosmer &al, 2013] refers to a type of categorical data where each observation can take one of two possible outcomes or categories. These outcomes are typically represented as 0 and 1, or as "success" and "failure," "yes" and "no," or any other appropriate labels.

Binary data is commonly encountered in various fields, including biology, social sciences, finance, and machine learning.

When analyzing binary data, various statistical methods can be employed. Some commonly used techniques include:

Proportions and percentages: Calculating the proportion or percentage of observations falling into each category.

Chi-square test: Assessing the independence or association between two categorical variables.

Logistic regression: Modeling the relationship between binary response variables and predictor variables.

Odds ratio: Measuring the odds of an event occurring in one category compared to another.

Binomial distribution: Modeling the probability of observing a specific number of successes in a fixed number of trials.

Bayesian analysis can also be applied to binary data, where prior distributions, likelihood functions, and posterior distributions are used to estimate parameters and make inferences.

**2.5 Vector data**:

Vector data [Bivand,& al,2013] is a type of spatial data representation that uses points, lines, and polygons to represent geographic features. It is commonly used in Geographic Information Systems (GIS) and spatial analysis. Vector data represents the real-world features by defining their geometry and attributes.

Here are the main components of vector data:

**Points**: Points represent individual locations or specific features with a single set of coordinates. They are often used to represent landmarks, cities, or sampling locations. Each point can have associated attributes such as a name, population, or temperature.

**Lines**: Lines represent linear features, such as roads, rivers, or pipelines. They are composed of a series of connected points. Lines can have attributes associated with them, such as road type, length, or speed limits.

**Polygons**: Polygons represent areas or regions. They are enclosed by a series of connected lines, forming a closed shape. Examples of polygons include countries, parks, or administrative boundaries. Polygons can also have associated attributes, such as area, population density, or land use.

Bayesian methods can be used with both binary and vector data. Here are some examples of how Bayesian methods can be applied to analyze binary and vector data:

3. **Bayesian methods with binary data:**

3.1 **Logistic regression**: Logistic regression is widely used for modeling binary data. In Bayesian logistic regression, a prior distribution is assigned to the coefficients of the model, and by applying Bayes' theorem, the posterior distribution is obtained. This posterior distribution enables predictions and hypothesis testing to be conducted.

The equation for logistic regression can be represented as (Hastie, 2009):

$$\text{logit}(p) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_n x_n \quad (2)$$

where:

logit (*p*): represents the logarithm of the odds of the binary outcome.

*p* : the probability of the binary outcome.

$\beta_0, \beta_1, \beta_2, ..., \beta_n$ : coefficients or parameters in logistic regression are associated with each independent variable.

$x_1, x_2, ..., x_n$ :are the values of the independent variables.

**3.2 Bayesian network:** A Bayesian network is a graphical model used to depict the probabilistic connections between a groups of variables. Bayesian networks can be used to model binary data by representing each binary variable as a node in the network and specifying conditional probabilities between the nodes[ Koller &al, 2009].

The equations of Bayesian networks rely on the principles of conditional probability and the chain rule of probability. These equations encompass various aspects of Bayesian networks, including the joint probability distribution, the conditional probability tables (CPTs), and the inference equations. Overall, a Bayesian network is a graphical model that effectively illustrates the probabilistic relationships among a set of variables [Jensen&al,2007].

1. The joint probability distribution: The joint probability distribution in a Bayesian network is determined by multiplying the conditional probabilities of each variable given its parents in the network.

Mathematically, we can represent it as:

$$P(X_1, X_2, ..., Xn) = P(X_1 \mid \text{Parents}(X_1)) * P(X_2 \mid \text{Parents}(X_2)) * ... * P(Xn \mid \text{Parents}(Xn)) \quad (3)$$

This equation represents the factorization property of Bayesian networks, where each variable's probability is conditioned on its parents.

2. Conditional probability tables (CPTs): CPTs are used to depict the conditional probabilities of each variable based on its parents. Each entry in the CPT specifies the probability distribution of a variable given the possible combinations of states of its parents. The CPTs provide the necessary information for updating probabilities and performing inference within the net.

3. Inference equation: inference equations are utilized to calculate the posterior probabilities of variables based on observed evidence. The primary objective of

inference is to determine the probability of a particular variable given evidence on other variables. This can be accomplished through Bayesian inference, which involves utilizing the joint probability distribution and the available evidence to calculate the desired posterior probability.

$$P(X \mid E) = \alpha * P(X, E) = \alpha * \Sigma_u P(X, u, E) \quad (4)$$

Where:

X: represents the target variable,

E: represents the evidence,

$\alpha$ : the normalization constant, and u represents the unobserved or hidden variables.

**4. Bayesian methods with vector data:**

**4.1 Bayesian linear regression**: Bayesian linear regression shares similarities with Bayesian logistic regression, but it is specifically employed to model continuous outcomes rather than binary outcomes. In this approach, a prior distribution is assigned to the coefficients of the linear regression model, and the posterior distribution is then derived using Bayes' theorem [Gelman,& al ,2013].

The equation for Bayesian linear regression can be described as [Bishop, 2006]:

$$y = X\beta + \varepsilon \quad (5)$$

since:

y: represents the observed values of the dependent variable.

X: The design matrix, also referred to as the matrix of independent variables, is structured in such a way that each row represents an observation and each column represents a predictor variable.

$\beta$ : The vector of unknown coefficients or parameters represents the association between the independent variables and the dependent variable.

$\varepsilon$ : The vector of random errors or noise term is supposedly to follow a normal distribution with a mean of zero and a constant variance.

At Bayesian linear regression, the prior distribution is specified for the coefficients $\beta$. Typically, a normal distribution or a multivariate normal distribution is used as the prior distribution for $\beta$.

**4.2 Gaussian process regression:** Gaussian process regression[Murphy, 2012] is a non-parametric method that can be used to model vector data. In this method, a prior distribution is placed on the function that maps the input vectors to the output vectors, and the posterior distribution is obtained using Bayes' theorem. Gaussian

process regression is often used in machine learning applications where the relationship between the input and output vectors is complex and unknown

In Gaussian Process Regression (GPR), the relationship between the input variables (often denoted as X) and the output variable (often denoted as y) is modeled as a draw from a Gaussian process. The basic equation for Gaussian process regression can be written as follows[Bishop, 2006]:

$$y = f(X) + \varepsilon \qquad (6)$$

since:

y : The array containing the observed values of the output variable.

X : The input variable matrix, where each row represents an observation and each column represents a feature or predictor variable.

f(X) : The underlying function, which is unknown, represents the relation between the input variables and the output variable. This function is modeled as a random draw from a Gaussian process.

ε : The vector of random errors or noise term is supposedly to follow a normal distribution with a mean of zero and a constant variance.

In order to completely define the Gaussian process, it is necessary to specify the mean function and covariance function (kernel). The mean function which represents the average or expected value of the underlying function f(X) when a specific input is given, while the covariance function determines the similarity or correlation between different input-output pairs.

Gaussian process regression aims to estimate the distribution of the output variable y for new, unseen inputs, based on a given set of observed input-output pairs (X, y). This is accomplished by calculating the posterior distribution of the Gaussian process, which is conditioned on the observed data. The posterior distribution combines the prior distribution of the Gaussian process with the likelihood of the observed data.

To make predictions for new inputs, Gaussian process regression uses the posterior distribution to estimate the mean and variance of the output variable at each input point. Where The mean prediction shows the expected value of the output variable, while the variance provides a measure of uncertainty or confidence in the prediction.

In summary, the equation for Gaussian process regression is a combination of the underlying function f(X) and a noise term ε, where f(X) is modeled as a draw from

a Gaussian process. The specific form of the Gaussian process and the choice of mean function and covariance function (kernel) depend on the problem at hand and the assumptions made about the underlying relationship between the input and output variables.

**In summary, Bayesian methods can be used with both binary and vector data. According to the type of problem and the nature of the data, the appropriate method is chosen.**

5. **Numerical cases**:

**5.1 Bayes with Binary Data**

Consider a dataset consisting of 1000 emails, where each email is classified as either spam (1) or not spam (0). The objective is to utilize Bayes' theorem to classify new emails as either spam or not spam based on their content. To accomplish this, we can represent each email as a binary vector, where each element of the vector corresponds to a word in the email. If a word appears in the email, the corresponding element in the vector is set to 1; otherwise, it is set to 0. Otherwise , We can then use the following steps:

Calculate the prior probabilities of spam and not spam emails in the dataset. Let's say that 100 of the 1000 emails are labeled as spam, Given the information, the initial probability of an email being classified as spam, denoted as P(spam), is 0.1, while the initial probability of an email being classified as not spam, denoted as P(not spam), is 0.9.

Calculate the likelihood probabilities of each word given spam and not spam emails. For example, suppose that the word "free" appears in 50 of the 100 spam emails, and in 5 of the 900 not spam emails. Then the likelihood probability of "free" given spam is P("free"|spam) = 0.5, and the likelihood probability of "free" given not spam is P("free"|not spam) = 0.005.

Given a new email, calculate the posterior probabilities of spam and not spam using Bayes' theorem. For example, suppose that the new email contains the words "free" and "buy". The posterior probability of an email being classified as spam can be computed using the following formula:

P(spam|"free", "buy") = P("free", "buy"|spam) * P(spam) / P("free", "buy")

= P("free"|spam) * P("buy"|spam) * P(spam) / P("free", "buy")

= 0.5 * 0.2 * 0.1 / P("free", "buy")

Similarly, the posterior probability of an email being classified as not spam can be determined using the following calculation:

P(not spam|"free", "buy") = P("free", "buy"|not spam) * P(not spam) / P("free", "buy")

= P("free"|not spam) * P("buy"|not spam) * P(not spam) / P("free", "buy")

= 0.005 * 0.01 * 0.9 / P("free", "buy")

The denominator P("free", "buy") is the probability of observing the words "free" and "buy" in any email, and can be calculated as follows:

P("free", "buy") = P("free", "buy"|spam) * P(spam) + P("free", "buy"|not spam) * P(not spam)

= P("free"|spam) * P("buy"|spam) * P(spam) + P("free"|not spam) * P("buy"|not spam) * P(not spam)

= 0.5 * 0.2 * 0.1 + 0.005 * 0.01 * 0.9

Once we have calculated the posterior probabilities of spam and not spam, we can classify the new email as spam if P(spam|"free", "buy") > P(not spam|"free", "buy"), and as not spam otherwise.

## 5.2 Bayes with Vector Data

Suppose we have a dataset of 1000 images, where each image is represented as a 28x28 pixel grayscale matrix, and labeled as either a digit from 0 to 9. We want to use Bayes' theorem to classify new images based on their pixel values. We can model each image as a 784-dimensional vector, where each element corresponds to a pixel value between 0 and 255. We can then use the next steps:

Determine the prior probabilities of each digit in the dataset .Let's say that there are 100 images of each digit, so the prior probability of each digit is P(digit) = 0.1.

Calculate the likelihood probabilities of each pixel value given each digit. For example, suppose that the pixel at location (i, j) has value 128 in 20 of the images of digit 3, and in 10 of the images of digit 8. Then the likelihood probability of pixel (i, j) having value 128 given digit 3 is P(pixel(i,j)=128|digit=3) = 0.2, and the likelihood probability of pixel (i, j) having value 128 given digit 8 is P(pixel(i,j)=128|digit=8) = 0.1.

Given a new image, calculate the posterior probabilities of each digit using Bayes' theorem. For example, suppose that the new image has pixel values asfollows:

pixel(1,1) = 100, pixel(1,2) = 200, ..., pixel(28,28) = 50

We can calculate the posterior probability of digit 3 as follows:

P(digit=3|image) = P(image|digit=3) * P(digit=3) / P(image)

where P(image|digit=3) is the likelihood probability of the image given digit 3, and P(image) is the probability of the image appearing in any digit category. We can calculate the likelihood probability as follows:

P(image|digit=3) = P(pixel(1,1)=100|digit=3) * P(pixel(1,2)=200|digit=3) * ... * P(pixel(28,28)=50|digit=3)

Similarly, we can calculate the posterior probability of each digit, and classify the new image as the digit with the highest posterior probability.

**In summary, Bayesian methods can be used with both binary and vector data. The selection of approach relies on the particular issue at hand and the characteristics of the data.**

### 6. Conclusion:

In conclusion, Bayesian methods are a powerful tool for classification problems with binary and vector data. These methods allow us to incorporate prior knowledge, handle missing data, and quantify uncertainty in our predictions. For binary data, Bayesian methods can be used to model the probability of each class given the observed binary values and make predictions for new observations. For vector data, Bayesian methods can be used to model the probability of each class given the observed vector values and make predictions for new observations.

The ability to incorporate prior knowledge is one of the primary benefits of Bayesian methods .This can be particularly useful in cases where we have some prior information about the problem, such as the expected distribution of the features or the prevalence of the classes. By incorporating this prior

## REFERENCES

1. Albert, J. H. (2009). Bayesian Computation with R. Springer.

2. Barber, D. (2012). Bayesian Reasoning and Machine Learning. Cambridge University Press.

3. Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.

4. Bivand, R. S., Pebesma, E. J., & Gómez-Rubio, V. (2013). Applied Spatial Data Analysis with R (2nd Edition). Springer.

5. Blei, D. M., Kucukelbir, A., & McAuliffe, J. D. (2017). Variational inference: A review for statisticians. Journal of the American Statistical Association, 112(518), 859-877.

6. Gelman, A., & Hill, J. (2007). Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press.Association 422, 669–679, 1993.

7. Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). Bayesian Data Analysis, Third Edition. Chapman and Hall/CRC.

8. Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence. Nature, 521(7553), 452-459.

9. Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition. Springer.

10. Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction (2nd Edition). Springer.

11. Hosmer Jr, D. W., Lemeshow, S., & Sturdivant, R. X. (2013). Applied Logistic Regression (3rd Edition). Wiley

12. Jensen, F. V., & Nielsen, T. D. (2007). Bayesian Networks and Decision Graphs. Springer.

13. Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. Science, 349(6245), 255-260.

14. Koller, D., & Friedman, N. (2009). Probabilistic Graphical Models: Principles and Techniques. MIT Press.

15. McElreath, R. (2020). Statistical Rethinking: A Bayesian Course with Examples in R and Stan (2nd Edition)

16. Murphy, K. P. (2012). Machine Learning: A Probabilistic Perspective. MIT Press.

## APPENDIX – CODES

Case 1: Bayes with Binary Data

Let's consider a dataset consisting of 1000 emails, with each email classified as either spam (1) or not spam (0). Our goal is to utilize Bayes' theorem to categorize new emails as either spam or not spam, based on their content. To achieve this, we can represent each email as a binary vector, where the i-th element of the vector is 1 if the i-th word appears in the email and 0 otherwise. Here's an example code:

Generate example data

```
set.seed(123)

n <- 1000

m <- 5000

X <- matrix(rbinom(n*m, 1, 0.1), nrow = n)

y <- rbinom(n, 1, 0.1)
```

Calculate prior probabilities

```
prior_spam <- sum(y == 1) / n

prior_not_spam <- 1 - prior_spam
```

Calculate likelihood probabilities

```
likelihood_spam <- apply(X[y == 1, ], 2, function(x) sum(x == 1) / sum(y == 1))

likelihood_not_spam <- apply(X[y == 0, ], 2, function(x) sum(x == 1) / sum(y == 0))
```

Define function to calculate posterior probabilities

```
calculate_posterior <- function(x} )

likelihood_spam_x <- likelihood_spam[x == 1]

likelihood_not_spam_x <- likelihood_not_spam[x == 1]

posterior_spam <- prod(likelihood_spam_x) * prior_spam

posterior_not_spam <- prod(likelihood_not_spam_x) * prior_not_spam

posterior_spam / (posterior_spam + posterior_not_spam)
```

Example usage of calculate_posterior

```
new_email <- c(rep(0, 2500), rep(1, 2500))
```

```
posterior_prob <- calculate_posterior(new_email)
```

```
if (posterior_prob > 0.5} )
```

```
cat("New email is spam\n("
```

```
 {else}
```

```
cat("New email is not spam\n)"
```

Case 2: Bayes with Vector Data

Suppose we have a dataset of 1000 images, where each image is represented as a 28x28 pixel grayscale matrix, and labeled as either a digit from 0 to 9. We want to use Bayes' theorem to classify new images based on their pixel values. We can model each image as a 784-dimensional vector, where each element corresponds to a pixel value between 0 and 255. Here's an example code:

Load MNIST dataset (requires 'mnist' package)

```
library(mnist)
```

```
train_images <- mnist$load_train()[[1]]
```

```
train_labels <- mnist$load_train()[[2]]
```

```
test_images <- mnist$load_test()[[1]]
```

```
test_labels <- mnist$load_test()[[2]][
```

Flatten images into vectors

```
train_images_vec <- apply(train_images, 3, function(x) as.vector(x))
```

```
test_images_vec <- apply(test_images, 3, function(x) as.vector(x))
```

Calculate prior probabilities

```
prior_digits <- table(train_labels) / length(train_labels)
```

Calculate likelihood probabilities

```
likelihood_digits <- lapply(0:9, function(digit} )
```

```
digit_images <- train_images_vec[train_labels = digit]
```

```
colMeans(digit_images)
```

Define function to calculate posterior probabilities

calculate_posterior <- function(x} )

likelihood_x <- unlist(lapply(0:9, function(digit} )

likelihood_digit_x <- likelihood_digits[[digit[[

prod(dnorm(x, mean = likelihood_digit_x, sd = 1, log = TRUE))

posterior_digits <- likelihood_x + log(prior_digits)

posterior_digits / sum(exp(posterior_digits))

Example usage of calculate_posterior

new_image <- test_images_vec[1 ],

posterior_probs <- calculate_posterior(new_image)

predicted_digit <- which.max(posterior_probs)

```
cat("Predicted digit:", predicted_digit, "\n")
cat("Posterior probabilities:", posterior_probs, "\n")
```