# STATISTICAL BAYESIAN ANALYSIS OF EXPERIMENTAL DATA.

**LABDAOUI AHLAM and MERABET HAYET**

Department of Mathematics, University Constantine 1
Route of Ain El Bey, 25000 Constantine, Algeria.

**Abstract :**

The Bayesian researcher should know the basic ideas underlying Bayesian methodology and the computational tools used in modern Bayesian econometrics. Some of the most important methods of posterior simulation are Monte Carlo integration, importance sampling, Gibbs sampling and the Metropolis- Hastings algorithm. The Bayesian should also be able to put the theory and computational tools together in the context of substantive empirical problems. We focus primarily on recent developments in Bayesian computation. Then we focus on particular models. Inevitably, we combine theory and computation in the context of particular models. Although we have tried to be reasonably complete in terms of covering the basic ideas of Bayesian theory and the computational tools most commonly used by the Bayesian, there is no way we can cover all the classes of models used in econometrics. We propose to the user of analysis of variance and linear regression model.

*Keywords*: Bayesian analysis, Markov Chain Monte Carlo Algorithms, regression models.

**Résumé**

Le chercheur bayésien devrait connaître les idées de base qui sous-tendent la méthodologie bayésienne et les outils informatiques utilisés dans l'économétrie bayésienne moderne. Certaines des méthodes les plus importantes de la simulation Monte Carlo postérieure sont l'intégration, l'échantillonnage d'importance, l'échantillonnage de Gibbs et l'algorithme de Metropolis-Hastings. Le filtre Bayésien doit également être en mesure de mettre la théorie et des outils informatiques ainsi que dans le contexte de fond les problèmes empiriques. Nous nous concentrons principalement sur les développements récents dans le calcul bayésien. Ensuite, nous nous concentrons sur des modèles particuliers. Inévitablement, nous combinons théorie et le calcul dans le contexte des modèles particuliers. Bien que nous ayons essayé d'être assez complet en termes de couverture des idées de base de la théorie bayésienne et les outils informatiques les plus couramment utilisées par le bayésien, il n'y a aucun moyen que nous pouvons couvrir tous les types de modèles utilisés en économétrie. Nous proposons à l'utilisateur de l'analyse de la variance et modèle de régression linéaire.

*Mots clés* : analyse bayésienne, les algorithmes de Monte Carlo par chaine de Markov, les modèles de régression.

**ملخص.**

يجب على الباحث البايزي معرفة المفاهيم الأساسية المعتمدة على الطريقة البايزية و أدوات الإعلام الآلي المستعملة في الإقتصاد البايزي المتطور, بعض الطرق المهمة في التقدير السالف لمونتكارلو هم: التكامل والعينات المهمة المتمثلة في عينة قيبس وخوارزمية ميتروبوليس هاستينك. المصفاة البايزية تؤخذ بعين الاعتبار في قياس و وضع النظرية وأدوات الإعلام الآلي وكذلك في كامل نص المشاكل التجريبية.ركزنا أساسا على التطور الجديد في الحساب البايزي. فيما بعد نركز على النماذج الخاصة, الحتمية وكذلك الحساب في نص النماذج الخاصة وحاولنا أيضا تكملة العبارات والأفكار الأساسية للنظرية البايزية وأدوات الإعلام الآلي الأكثر استعمالا من طرف البايزي, لا توجد أي وسيلة تمكننا من حجب جميع أنواع النماذج المستعملة في الاقتصاد لذا نقترح على المستعمل تحليل التت تت ونموذج التراجع الخطي.

*الكلمات المفتاحية* : التحليل البايزي , خوارزمية مونتكارلوبسلسلة ماركوف, نماذج التراجع.

# Introduction :

Regression is by far the larger the field of statistics, both theoretical and applied. This is the preferred method of econometrics, and the practice of social science modeled on econometrics, "econometric model" has come to mean any regression model, even without reference to economic problems.

The framework model of regression is defined by a variable to predict (or "dependent", dedicated notation y), and a variable (simple regression) and multivariate (multiple regression) known predictor variables (or "independent"). Regression is to construct a variable regressed $\widehat{y}$ combination of predictor variables as close as possible (in a sense to be specified) of the dependent variable.

Procedures classical linear regression, applicable to numeric variables, recently came to enlist the logistic regression and its variants for the variables categorized. Considerations of this module focus on linear regression, shall apply (mutatis mutandis) to various forms of regression.

Statistical experimental data, regression can be considered as a special case of the analysis of variance, in the case of digital independent variables. For observational data, new problems arise, related to the fact that in general the predictor variables are not statistically independent. It is these problems that have focused my recent work.

In the Bayesian framework, there is no fundamental difference between the observation and the parameter of a statistical model, both of which are considered variable quantities, so if we denote by x the given bill sampling $f(x \mid \theta)$, and $\theta$ the model parameters considered (plus possibly latent variables) of prior formal inference requires updating of the conditional distribution $f(\theta \mid x)$ parameter. Determining $\pi(\theta)$ and $f(x \mid \theta)$ gives $f(x, \theta)$ by

$$f(x, \theta) = f(x \mid \theta) * \pi(\theta) \qquad (1)$$

After observing x, we can use Bayes' theorem to determine the distribution of $\theta$ conditional on the data (or the posterior) (see [2]).

$$\pi(\theta \mid x) = \frac{f(x \mid \theta) * \pi(\theta)}{\int f(x \mid \theta) * \pi(\theta) d(\theta)} \qquad (2)$$

For the Bayesian approach, all the features of the posterior distribution are important for inference: time, quantile, etc ... These quantities can often be expressed in terms of conditional expectation of a function of $\theta$ with respect to the law post (see [2])

$$E(h(\theta) \mid x) = \frac{\int h(\theta) f(x \mid \theta) * \pi(\theta) d(\theta)}{\int f(x \mid \theta) * \pi(\theta) d\theta} \qquad (3)$$

We can calculate the posterior distribution directly in the simple case or calculation is made by MCMC simulation where the integral calculation is very complex.

In our work we first present the regression model and the simple and multiple logistic model then we set the conditions for the use of algorithms Monte Carlo Markov Chain (MCMC) then we introduce some MCMC algorithms, in particular the Metropolis-Hastings algorithm and the Gibbs sampling method. Finally, we present the numerical results and their interpretations.

We used the software WinBUGS to estimate the parameters, and interpret the results of actual data, WinBUGS (the MS Windows operating system version of BUGS: Bayesian
Analysis Using Gibbs Sampling) is a versatile package that has been designed to carry out Markov chain Monte Carlo (MCMC) computations for a wide variety of Bayesian models(see[6]).

## 2 Methodology

## 2.1 Regression models

## 2.1.1. Linear regression model

Regression is for a type of problem where two continuous quantitative variables X and Y have a role asymmetrical variable Y depends on the variable X. The connection between the dependent variable Y and the independent variable X can be modeled as a function of Y = α + β X, (see [3])

Y: dependent variable (explained)

X: independent variable (predictor)

α: intercept (value of Y for x = 0)

β: slope (average variation of the value of Y for a one-unit increase of α et β can be calculated by :

$$\beta = \text{regression coefficient} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sum(X - \bar{X})^2} \qquad (4)$$

$$\alpha = \text{Intercept} = \bar{Y} - \beta \bar{X} \qquad (5)$$

r = correlation coefficient = is another important determinant and looks a lot like β.

$$r=\frac{\sum(X-\bar{X})(Y-\bar{Y})}{\sqrt{\sum(X-\bar{X})^2\sum(Y-\bar{Y})^2}} \qquad (6)$$

r = measure for the strength of association between Y and X-data. The stronger the association, the better Y predicts X.

### 2.1.2 Multiple linear models

The multiple regression model is a generalization of the regression model Simple when the explanatory variables are finite in number. The connection between the dependent variable Y and the independent variables X1 and X2 can be modeled as a function of

Y= α + β₁* X₁ + β₂* X₂.

A linear regression model is defined by an equation of the form:

$$Y_{n\times1} = X_{n\times p}\beta_{p\times1} + \varepsilon_{n\times1} \qquad (7)$$

Y : is an n-dimensional random vector.

X: is a matrix of size n × p known design matrix called experience.

β: is the p-dimensional vector of unknown model parameters

ε: the vector is centered, n-dimensional errors.

### 2.1.3 Logistic model

A standard qualitative regression and logistic regression model or logit model, where the conditional distribution of y is z∈$R^p$ explanatory variables, (see [2]):

$$P(y=1) = 1 - P(y=0) = \frac{\exp(z^t\gamma)}{1+\exp(z^t\gamma)} \qquad (8)$$

Consider the particular case where z= (1,$^x$) and $^Y$= (α, β) random variables $y_i$ values in {0,1} are associated with explanatory variables were modeled using a Bernoulli conditional probability

$$y_i\backslash x_i \sim B\left(\frac{\exp(\alpha+\beta x_i)}{1+\exp(\alpha+\beta x_i)}\right) \qquad (9)$$

Assume that our parameters follow a priori law unsuitable $\pi$(α, β) = 1. The likelihood of our model for a sample ($^{y_1}$, $^{x_1}$),…,($^{y_n}$,$^{x_n}$),is equal to

$$f(y_1,\dots,y_n\backslash x_1,\dots\dots,x_n,\alpha,\beta) = \prod_{i=1}^n \frac{\exp\{(\alpha+\beta x_i)y_i\}}{1+\exp(\alpha+\beta x_i)} \qquad (10)$$

The posterior distribution of (α, β) is then deduced by formal application of Bayes Theorem, see [10]

$$\alpha \prod_{i=1}^n \frac{\exp\{(\alpha+\beta x_i)y_i\}}{1+\exp(\alpha+\beta x_i)} = \frac{exp\{\sum_{i=1}^n(\alpha+\beta x_i)y_i\}}{\prod_{i=1}^n 1+\exp(\alpha+\beta x_i)} \qquad (11)$$

### 2.2. MCMC methods

The Monte Carlo Markov Chain (Monte Carlo Markov Chains in English or MCMC) is used when interest law cannot be simulated directly by the usual methods and / or when its density is known to a normalization constant fields. (see [4])

### 2.2.1. Metropolis-Hasting algorithm

The Metropolis-Hastings algorithm based on the use of a conditional density measurement $q(y|x)$ with respect to the dominant model li. It cannot be put into practice if $q(.|x)$ is simulated quickly and is available either analytically for a constant independent of either symmetrical, that is to say as $q(y|x) = q(x|y)$. The Metropolis-Hastings algorithm (see [5]) associated with the objective law $\pi$ and the conditional $q$ produces a Markov chain $x^{(t)}$ based on the following transition:

Initialization: X₀
At each step k ≥ 0:
  • Simulate a value
  • Simulate a value.
  • Ask

$$X_{k+1}\begin{cases} y_k \ if \ u_k \leq \rho(x_k,y_k) \\ X_k \ else, \end{cases}$$

$$\rho(x_k,y_k) = \min\left\{1,\frac{\pi(y_k)q(x_k|y_k)}{\pi(x_k)q(y_k|x_k)}\right\}.$$

Or

The law $q$ is called the law of instrumental or proposal. This algorithm accepts systematically simulations $y_t$ such that the ratio $\left(\pi(y_t)\big|q(y_t|x^{(t)})\right)$ is greater than the previous value $\left(\pi((x^{(t)}))\big|q(x^{(t)}|y_t)\right)$. It is only in

the symmetric case that acceptance is governed by the report $\pi(y_t)/\pi(x_t)$.

### 2.2.2 The Gibbs sampling

The Gibbs sampling algorithm is a simulation of a law $\pi(x)$ such that:

x admits a decomposition of the form

$$x = (x_1, \ldots, x_n),$$

The conditional law

$$\pi_i \ (.|(x_1, \ldots, x_{x-1}, x_{x+1}, \ldots, x_n))$$ are easily

simulated (see[9]).

Example: $(X, Y) \sim N_{(0, \Sigma)}$, with $\Sigma = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$.

$$X \mid Y = y \sim N(\rho x, 1 - \rho^2)$$

Principle of the algorithm: Updating "component by component".

At each step $k \geq 0$ :

- simulate $X_1^{k+1} \sim \pi_1\left(. \left| X_2^k, \ldots \ldots \ldots, X_n^k \right.\right)$

- simulate
$$X_i^{k+1} \sim \pi_i\left(. \left| X_1^k, \ldots \ldots, X_{i-1}^{k+1}, X_{i-1}^k, \ldots \ldots, X_n^k \right.\right)$$

- simulate
$$X_n^{k+1} \sim \pi_1\left(. \left| X_1^{k+1}, \ldots \ldots \ldots, X_{n-1}^{k+1} \right.\right)$$

## 3 - Applications

### 3.1 - Example with WinBUGS: linear model "calculated α and β"

Table 1 gives the real data of a crossover study comparing a new laxative versus a standard laxative, bisacodyl. Days with stool are used as primary endpoint. The table shows that the new drug is more efficacious than bisacodyl (see [11]).

| Patient no | Y-variables | X-variables |
|---|---|---|
| 1 | 24 | 8 |
| 2 | 30 | 13 |
| 3 | 25 | 15 |
| 4 | 35 | 10 |
| 5 | 39 | 9 |
| 6 | 30 | 10 |
| 7 | 27 | 8 |
| 8 | 14 | 5 |
| 9 | 39 | 13 |
| 10 | 42 | 15 |
| 11 | 41 | 11 |
| 12 | 38 | 11 |
| 13 | 39 | 12 |
| 14 | 37 | 10 |
| 15 | 47 | 18 |
| 16 | 30 | 13 |
| 17 | 36 | 12 |
| 18 | 12 | 4 |
| 19 | 26 | 10 |
| 20 | 20 | 8 |
| 21 | 43 | 16 |
| 22 | 31 | 15 |
| 23 | 40 | 14 |
| 24 | 31 | 7 |
| 25 | 36 | 12 |
| 26 | 21 | 6 |
| 27 | 44 | 19 |
| 28 | 11 | 5 |
| 29 | 27 | 8 |
| 30 | 24 | 9 |
| 31 | 40 | 15 |
| 32 | 32 | 7 |
| 33 | 10 | 6 |
| 34 | 37 | 14 |
| 35 | 19 | 7 |

*Table 1: Example of a crossover trial comparing efficacy of a new laxative versus bisacodyl*

*Model with software WinBUGS

    Y-variables: new treatment (days with stool).
    X-variables: bisacodyl (days of stool).

$Y \sim N(mu_i, tau)$
$mu_i = \alpha + \beta * X_i$

The model is:

```
model
{
  for(i in 1 : 35) {
  y[i] ~ dnorm(mu[i], tau)
  mu[i] <- alpha + beta * X[i]
}
alpha ~ dnorm(0, 1.0E-6)
beta ~ dnorm(0, 1.0E-6)
tau ~ dgamma(1.0E-3, 1.0E-3)
sigma <- 1/sqrt(tau)
}
```

We then proceed to estimate, this time on two channels, with 110 000 iterations (1000 enough) each, keeping an iteration of 150. The parameters of the line are

estimated, α = 8.669 with a standard deviation of 3.236 and β= 2.062 with a standard deviation of 0.2854. WinBUGS outputs are as follows:

| node | mean | sd | MC error |
|------|------|-----|----------|
| alpha | 8.669 | 3.236 | 0.02289 |
| beta | 2.062 | 0.2854 | 0.002044 |
| tau | 0.02641 | 0.006571 | 4.7E-5 |

| 2.5% | median | 97.5% | start | sample |
|------|--------|-------|-------|--------|
| 2.308 | 8.658 | 15.03 | 1 | 22000 |
| 1.505 | 2.061 | 2.624 | 1 | 22000 |
| 0.01533 | 0.02583 | 0.04088 | 1 | 22000 |

We now presenting a graphical representation of the parameters alpha and beta, of Kernel density in *fig.1*, quantiles in *fig.2* and the auto correlation function in *fig.3*



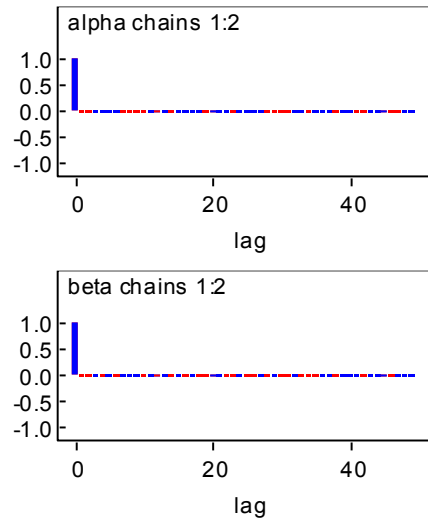*Figure 1 :kernel density*



*Figure 2 : Quantiles*



*Figure 3 : Autocorrelation function*

## 3.2 Example with WinBUGS: multiple linear model "calculated α, β₁ and β₂"

We may be Interested to know if age is an independent contributor to the effect of the new laxative. That purpose for a simple regression equation has to be extended as follows $Y = \alpha + \beta_1 X_1 + \beta_2 X_2$, two partial regression coefficients are Called. Just like a simple linear regression, multiple linear regression can give us the best fit for the data given, although it is hard to display the correlations in a figure. Table 2 gives the data from Table 1 extended by the variable age (see [11]).

| Patient no | Y-variables | X₁-variables | X₂-variables |
|-----------|-------------|-------------|-------------|
| 1 | 24 | 8 | 25 |
| 2 | 30 | 13 | 30 |
| 3 | 25 | 15 | 25 |
| 4 | 35 | 10 | 31 |
| 5 | 39 | 9 | 36 |
| 6 | 30 | 10 | 33 |
| 7 | 27 | 8 | 22 |
| 8 | 14 | 5 | 18 |
| 9 | 39 | 13 | 14 |
| 10 | 42 | 15 | 30 |
| 11 | 41 | 11 | 36 |
| 12 | 38 | 11 | 30 |
| 13 | 39 | 12 | 27 |
| 14 | 37 | 10 | 38 |
| 15 | 47 | 18 | 40 |
| 16 | 30 | 13 | 31 |
| 17 | 36 | 12 | 25 |
| 18 | 12 | 4 | 24 |
| 19 | 26 | 10 | 27 |
| 20 | 20 | 8 | 20 |

| 21 | 43 | 16 | 35 |
|----|----|----|----|
| 22 | 31 | 15 | 29 |
| 23 | 40 | 14 | 32 |
| 24 | 31 | 7 | 30 |
| 25 | 36 | 12 | 40 |
| 26 | 21 | 6 | 31 |
| 27 | 44 | 19 | 41 |
| 28 | 11 | 5 | 26 |
| 29 | 27 | 8 | 24 |
| 30 | 24 | 9 | 30 |
| 31 | 40 | 15 | 20 |
| 32 | 32 | 7 | 31 |
| 33 | 10 | 6 | 29 |
| 34 | 37 | 14 | 43 |
| 35 | 19 | 7 | 30 |

| node | mean | sd | MC error | 2.5% | median | 97.5% | start | sample |
|------|------|-----|----------|------|--------|-------|-------|--------|
| alpha | 2.332 | 4.985 | 0.03595 | -7.53 | 2.353 | 12.14 | 1 | 22000 |
| beta1 | 1.876 | 0.3003 | 0.002056 | 1.282 | 1.875 | 2.472 | 1 | 22000 |
| beta2 | 0.2827 | 0.1719 | 0.001191 | -0.05832 | 0.2827 | 0.6232 | 1 | 22000 |
| tau | 0.02786 | 0.006994 | 5.042E-5 | 0.01588 | 0.02728 | 0.04326 | 1 | 22000 |

Table 2: *Example of a crossover trial comparing efficacy of a new laxative versus bisacodyl*

*Model with software WinBUGS

Y-variables: new treatment (days with stool).
X1-variables: bisacodyl (days of stool).
X2-variables: age (years).

$Y \sim N (mu, tau)$
$mu = \alpha + \beta_1 * X_1 + \beta_2 * X_2$

The model is:

model

```
    {
     for(i in 1 : 35) {
       y[i] ~ dnorm(mu[i], tau)
       mu[i] <- alpha + beta1* X1[i]  +
beta2* X2[i]
     }
    alpha ~ dnorm(0, 1.0E-6)
    beta1 ~ dnorm(0, 1.0E-6)
     beta2 ~ dnorm(0, 1.0E-6)
    tau ~ dgamma(1.0E-3, 1.0E-3)
    sigma <- 1/sqrt(tau)
    }
```

We now presenting a graphical representation of the parameters alpha and beta (1), beta (2) of Kernel density in *fig.4*, quantiles, *fig.5* and the graphical of auto correlation in *fig.6*:
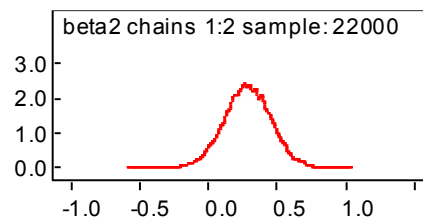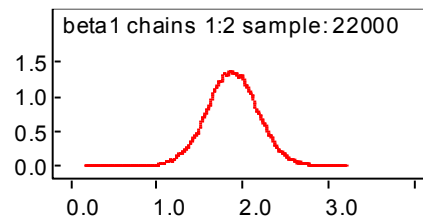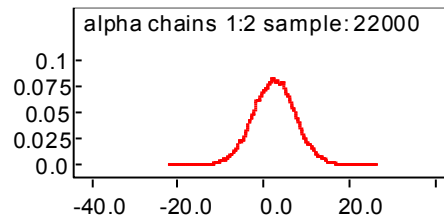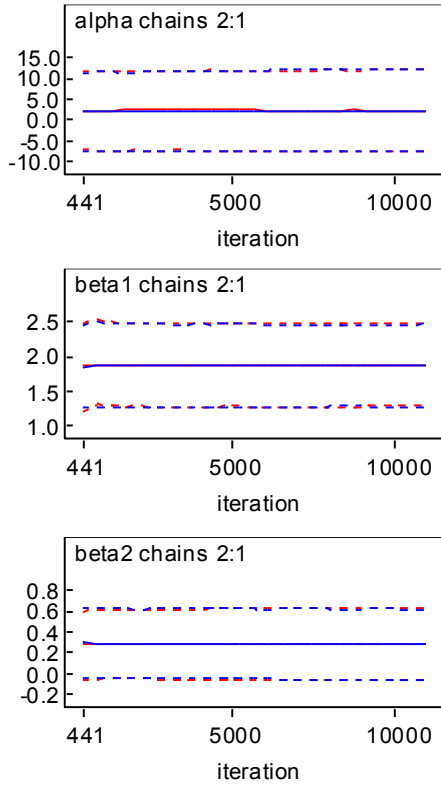


*Figure4: Kernel density*

We then proceed to estimate, this time on two channels, with 110 000 iterations (1000 enough) each, keeping an iteration of 150. The parameters of the line are estimated, $\alpha = 2.332$ with a standard deviation of 4.985 and $\beta_1 = 1.876$ with a standard deviation of 0.3003, $\beta_2 = 0.282$ with a standard deviation of 0.171.
WinBUGS outputs are as follows:
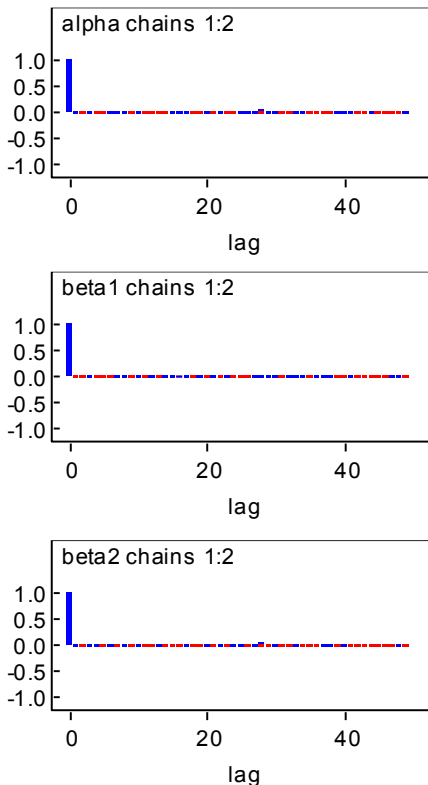
*Figure5: quantiles*



*Figure6: Autocorrelation function*

## 3.3 Example with WinBUGS: logistic model

Our study is based on a comparison of an antiseptic cream and Placebo; as the endpoint is cure an infection. We seek to estimate the effect of the cream versus placebo, the following table gives the answer 8 centers that we have considered, see [1]:

| Centre | traitement | Réponse | |
|--------|-----------|---------|-------|
| | | Succès | Echec |
| 1 | Crème | 11 | 25 |
| | Placebo | 10 | 27 |
| 2 | Crème | 16 | 4 |
| | Placebo | 22 | 10 |
| 3 | Crème | 14 | 5 |
| | Placebo | 7 | 12 |
| 4 | Crème | 2 | 14 |
| | Placebo | 1 | 16 |
| 5 | Crème | 6 | 11 |
| | Placebo | 0 | 12 |
| 6 | Crème | 1 | 10 |
| | Placebo | 0 | 10 |
| 7 | Crème | 1 | 4 |
| | Placebo | 1 | 8 |
| 8 | Crème | 4 | 2 |
| | Placebo | 6 | 1 |

Table3: processed data

$$* \, r_i^T \rightarrow \text{Binomial} \, (p_i^T, n_i^T),$$
$$* \, \text{logit} \, (P_i^p) = \alpha - \beta/2 + u_i,$$
$$* \, \text{logit} \, (P_i^c) = \alpha + \beta/2 + u_i,$$
$$* \, u_i \rightarrow \text{Normal} \, (0, \sigma_u^2)$$

The model is:

```
model
{
for(i in 1 : 8) {
rp[i] ~ dbin(pp[i], np[i])
rc[i] ~ dbin(pc[i], nc[i])
logit(pp[i]) <- alpha - beta / 2 + u[i]
logit(pc[i]) <- alpha + beta / 2 + u[i]
u[i] ~ dnorm(0.0, tau)
}
alpha ~ dnorm(0.0, 1.0E-6)
beta ~ dnorm(0.0, 1.0E-6)
tau ~ dgamma(0.1, 0.1)
sigma <- 1/ sqrt(tau)
OR <- exp(beta)
}
```

We then proceed to estimate, this time on three channels, with 110 000 iterations (1000 enough) each, keeping an iteration of 150. The (assumed homogeneous) cream is estimated at 0.757, with a standard deviation of 0.304. WinBUGS outputs are as follows:

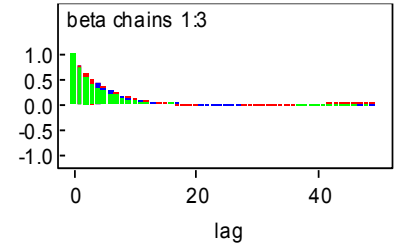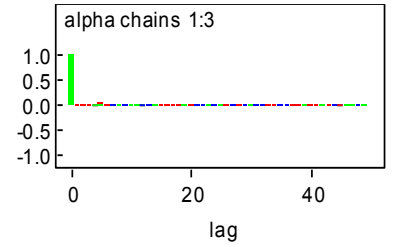| node | mean | sd | MC error | 2.5% | median | 97.5% | start | sample |
|------|------|-----|----------|------|--------|-------|-------|--------|
| OR | 2.235 | 0.7052 | 0.009858 | 1.191 | 2.123 | 3.88 | 1 | 3000 |
| alpha | -0.8362 | 0.757 | 0.05594 | -2.37 | -0.859 | 0.8177 | 1 | 3000 |
| beta | 0.7575 | 0.3043 | 0.004185 | 0.1745 | 0.7528 | 1.356 | 1 | 3000 |
| tau | 0.4806 | 0.7091 | 0.01471 | 0.0881 | 0.3918 | 1.323 | 1 | 3000 |
| u[1] | -0.1132 | 0.7887 | 0.05607 | -1.787 | -0.09136 | 1.514 | 1 | 3000 |
| u[2] | 1.893 | 0.8025 | 0.05561 | 0.2114 | 1.89 | 3.542 | 1 | 3000 |
| u[3] | 1.01 | 0.8064 | 0.05598 | -0.7026 | 1.008 | 2.671 | 1 | 3000 |
| u[4] | -1.427 | 0.9 | 0.05286 | -3.341 | -1.3820. | 2316 | 1 | 3000 |
| u[5] | -0.6119 | 0.8564 | 0.05491 | -2.436 | -0.5866 | 1.111 | 1 | 3000 |
| u[6] | -1.871 | 1.044 | 0.04911 | -4.184 | -1.764 | -0.05686 | 1 | 3000 |
| u[7] | -0.849 | 0.9642 | 0.05094 | -2.914 | -0.7977 | 0.9901 | 1 | 3000 |
| u[8] | 1.856 | 0.9282 | 0.05278 | 0.1017 | 1.806 | 3.815 | 1 | 3000 |

We now presenting a graphical representation of the parameters alpha and beta, of Kernel density in *fig.7*, quantiles *fig.8* and finally the auto correlation function in *fig.9*
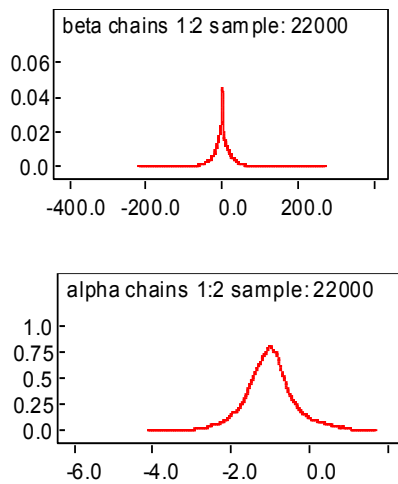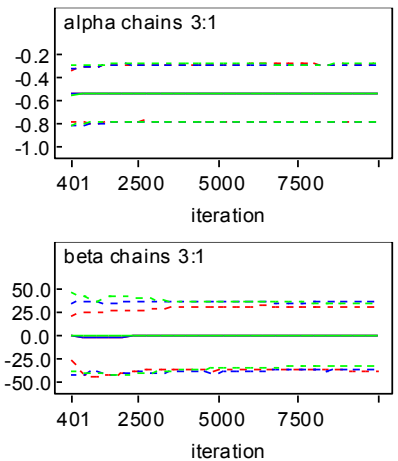


Figure 7 : Kernel density



Figure 8: quantiles



Figure 9: Autocorrelation function

## 4 - Discussion

* In the linear regression model the regression line is Y = 2.065+X 8.646.
The slope is 2.065 and directed the original is 8.646, and if we have x = 1 → y = 10 therefore
the new treatment is better than the standard treatment.

** The regression line in the model of multiple linear regression is Y = 2332 + 0282 + X2 1.876X1
we add the parameter age or not the new treatment is the best.

*** As gold is greater than 1 and the confidence interval between 3.88 and 1191 at 97.5% it is said that our anti septic cream is effective.

## 5 – Conclusion :

One of the merits of our work is to have shown using experimental data of clinical trials that can be modeled in a natural way and draw appropriate inferences, namely estimating parameters in regression models: model simple and multiple linear and logit model using Monte Carlo methods for Markov Chain (MCMC) especially as computer performance, made feasible processes effective simulations and the availability of computer programs has facilitated the calculation of posterior probabilities, which were previously daunting complexity.

## 6 – References :

[1] Agresti A. *Categorical Data Analysis*. (2002).

[2] Anas Altaleb, Christian P. Robert, *Analyse bayésienne du modèle logit : algorithme par tranches ou Metropolis-Hastings?*, revue de statistique appliquée, tome 49, n°4 (2001), p. 53-70.

[3] Christian P. Robert, Jean-Michel Marin, Bayesian Core: A Practical Approach to Computational Bayesian Statistics

[4] Christian P, Robert and George Casella, *Monte Carlo Statistical Methods*, Springer, (2004).

[5] C. Robert et G. Casella, *Monte Carlo Statistical Methods*, Springer, 2nd edition, (2004).

[6] David J. Lunn, Andrew Thomaa, Nicky Best and David Spiegelhalter *WinBUGS – A Bayesian modeling framework: Concepts, structure, and extensibility, Statistics and Computing* (2000) **10,** 325–337.

[7] Éric. Parent. Jacques Bernier*,   Le raisonnement bayésien*, Springer-Verlag France, Paris, (2007).

[8] Lionel Riou França, *statistique bayésienne*, INSERM U669, Mai 2009.

[9] Robert, C.P. and Casella, G. *Monte Carlo Statistical Methods*. New York: Springer Verlag (1999).

[10] Robert, C.P. *L'analyse statistique bayésienne Economica*, Paris (1992).

[11] Ton J Cleophas, Aeilko H Zwinderman, Toine F Cleophas, *Statistics Applied to Clinical Trials* (2006).