

Get started

Open in app



Follow

625K Followers



You have **2** free member-only stories left this month. [Sign up for Medium and get an extra one](#)

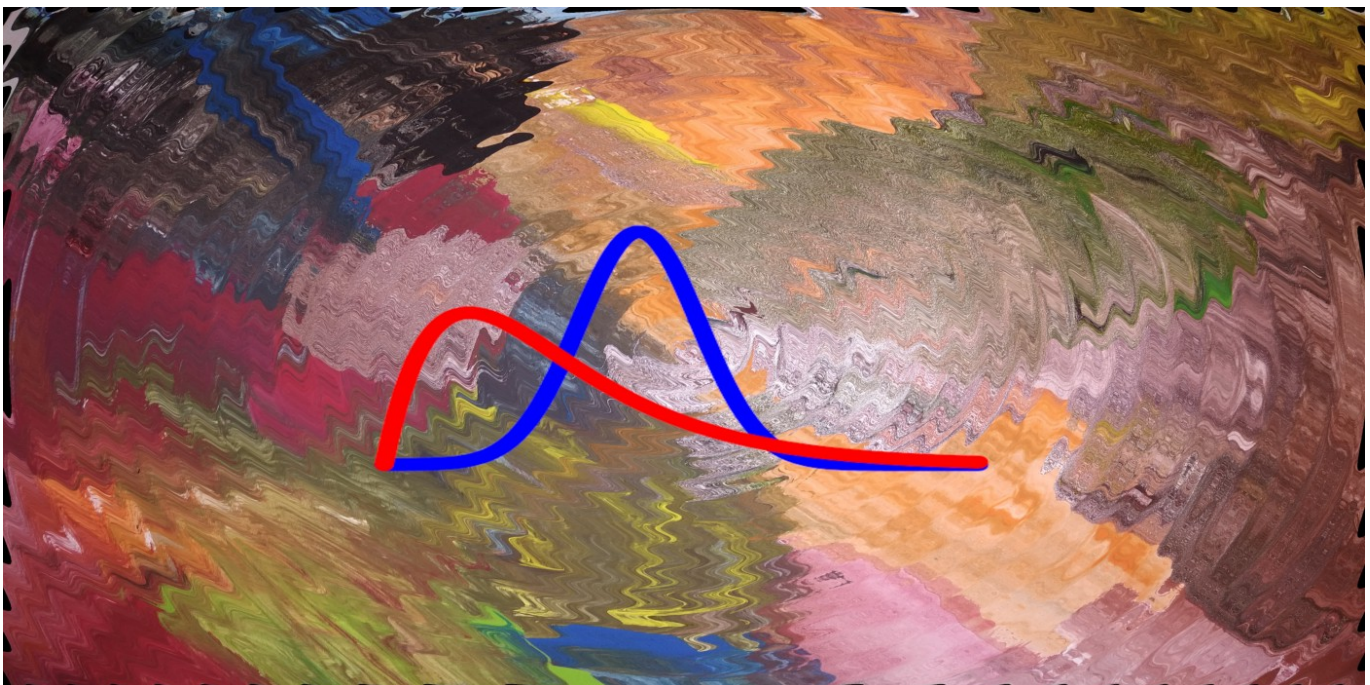
# Introduction to Bayesian Logistic Regression

A practical demonstration of the Bayesian approach to classification using Python and PyJAGS.



Michel Kana, Ph.D Feb 7, 2020 · 8 min read ★

This article introduces everything you need in order to take off with Bayesian data analysis. We provide a step-by-step guide on how to fit a Bayesian logistic model to data using Python. You will be able to understand Bayesian fundamentals for classification without dealing with math.



[Get started](#)[Open in app](#)

Let's review the concepts underlying Bayesian statistical analysis by walking through a simple classification model.

## The data

The data come from the [1988 Bangladesh Fertility Survey](#), where 1934 observations were taken from women in urban and rural areas. The authors of the dataset, Mn and Cleland aimed to determine trends and causes of fertility as well as differences in fertility and child mortality.

We will use the data in order to train a Bayesian logistic regression model that can predict if a given woman uses contraception.

The dataset is well suited to Bayesian logistic regression because being able to quantify uncertainty when analyzing fertility is the major component of population dynamics that decide the size, structure, and composition of populations ([source 1](#), [source 2](#)).

A copy of the raw data can be found [here](#). There are four attributes for each woman, along with a label indicating if she uses contraceptives. The attributes include:

- **district**: identifying code for the district the woman lives in,
- **urban**: type of region of residence,
- **living.children**: number of living children,
- **age-mean**: age of the woman (in years, centered around mean).

```
1 df_contraceptives = pd.read_csv('dataset_bang_contraceptive.csv')
2 df_contraceptives.head()
```

bayes\_logreg\_data.py hosted with ❤ by GitHub

[view raw](#)

district	urban	living.children	age_mean	contraceptive_use
35	0	4	2.4400	0
22	0	2	-1.5599	1
29	0	2	-8.5599	1
5	0	3	-4.5599	1

[Get started](#)[Open in app](#)

The women are grouped into 60 districts.

```
1 districts = df_contraceptives.sort_values(by="district").district.unique()
2 nb_districts = len(districts)
3 districts
```

bayes\_logreg\_districts1.py hosted with ❤ by GitHub

[view raw](#)

```
array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
        18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
        35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
        52, 53, 55, 56, 57, 58, 59, 60, 61])
```

A little bit of pre-processing is required. We map the district number 61 to the number 54 so that the districts are in order, as you can see below.

```
1 df_contraceptives.replace({'district': {61: 54}}, inplace=True)
2 df_contraceptives.sort_values(by="district").district.unique()
```

bayes\_logreg\_districts2.py hosted with ❤ by GitHub

[view raw](#)

```
array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
        18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
        35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
        52, 53, 54, 55, 56, 57, 58, 59, 60])
```

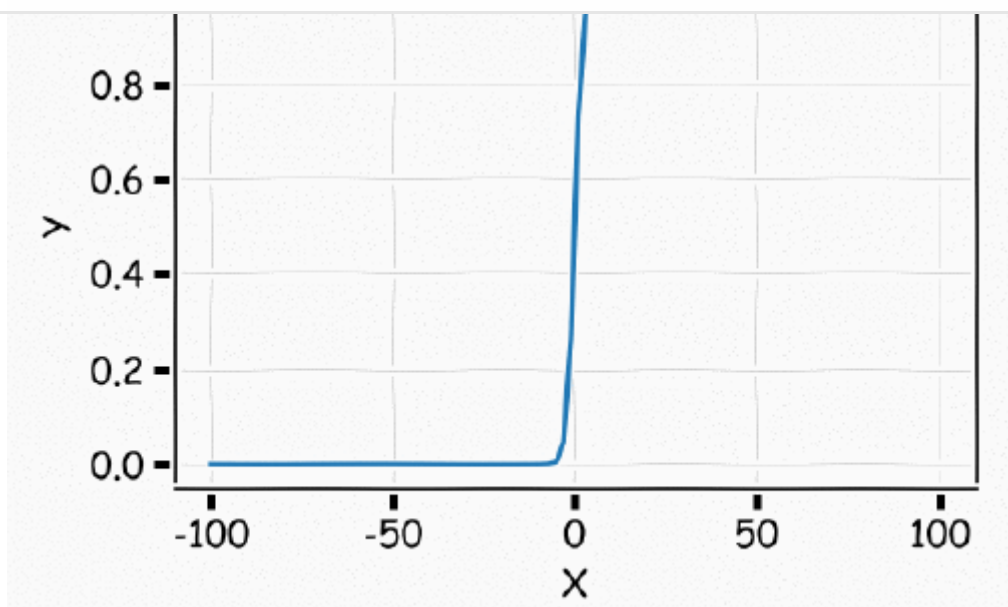
## Logistic regression

Predicting whether or not a given woman uses contraceptives is an example of binary classification problem. If we denote attributes of the woman by  $X$  and the outcome by  $Y$ , then the likelihood of using contraceptives,  $P(Y=1)$ , would follow the logistic function below.

$$P(Y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}}$$

Get started

Open in app



source: P. Protopapas, Harvard

The outcome depends on two parameters  $\beta_0$  (intercept) and  $\beta_1$  (slope). The intercept shifts the curve right or left, while the slope controls how steep the S-shaped curve is. If  $\beta_1$  is positive, then the predicted  $P(Y=1)$  goes from zero for small values of  $X$  to one for large values of  $X$  and if  $\beta_1$  is negative, then has the  $P(Y=1)$  opposite association.

For example, if  $\beta_1$  is positive for the *age* predictor, it means that older women are more likely to use contraceptives than younger ones. The cut-off threshold where *age* will start increasing the probability of contraceptive usage would be governed by  $\beta_0$ . Small values of  $\beta_0$  would indicate that contraceptives usage is a widespread practice in the population.

Understanding the logistic function is important for motivating the Bayesian approach. You can learn more about classical logistic regression in our article below.

### Why Deep Learning Works: solving a farmer's problem

In the beginning was the neuron: a gentle review of gradient descent, backpropagation, regression, autoencoder, CNNs...

[towardsdatascience.com](https://towardsdatascience.com)

[Get started](#)[Open in app](#)

and  $\beta_1$ . However, we will not know how confident we are about these estimates.

Moreover,  $\beta_0$  would normally differ from one district to another in real life.  $\beta_0$  could be small in some districts, but large in others. A classical logistic regression model would still provide a single value for all regions, which could lead to wrong conclusions.

In one of our past articles, we highlighted issues with uncertainty in machine learning and introduced the essential characteristics of Bayesian methods. We gently explained the explicit use of probability for quantifying uncertainty in inferences based on statistical data analysis.

### Bayesian nightmare. Solved!

Gentle introduction to Bayesian data analysis by examples and code in Python PyMC3.

[towardsdatascience.com](https://towardsdatascience.com)

Following Bayes, ideally, we want to take prior information into consideration when building our model for predicting contraceptive usage. If we only use the logistic equation, there is no direct way to include our **prior belief** about those parameters we are estimating.

*The Bayesian approach allows us to make a prior good guess of the intercept and slope, based on our real-life domain knowledge and common sense. We can say for example that, from experience, the intercept is drawn from a normal distribution with mean  $\mu=2$  and standard deviation  $\sigma=1$ . We can also postulate that the slope for the predictors *urban*, *living.children*, and *age-mean* are 4, -3 and -2 respectively. This prior belief is summarized as follows.*

$$\mu_{\beta_0} = 2$$

$$\sigma_{\beta_0}^2 = 1$$

$$\beta_0 \sim N(\mu_{\beta_0}, \sigma_{\beta_0}^2)$$

Get started

Open in app



$$\beta_3 = -2$$

Given the above distribution, which describes our prior belief, we can generate simulated data using a so-called **generative model**, as depicted in the image below.

$$Y_{ij} \sim \text{Bernoulli}(p_{ij})$$

$$\text{logit } p_{ij} = \beta_{0j} + \beta_1 \times \text{urban} + \beta_2 \times \text{living.children} + \beta_3 \times \text{age-mean},$$

where  $Y_{ij}$  is 1 if woman  $i$  in district  $j$  uses contraceptives, and 0 otherwise, and where  $i \in \{1, \dots, N\}$  and  $j \in \{1, \dots, J\}$ .  $N$  is the number of observations in the data, and  $J$  is the number of districts.

In this case, the generative model is a Bernoulli distribution with parameters  $p$ , which is the probability of a woman using contraceptives.  $p$  is the logistic function introduced in the previous section.

In order to get a better grasp of the concept of generative model, let's simulate binary response data  $Y$ . We do this by using prior parameter values and data.

```

1  beta_0_mu = 2
2  beta_0_sigma = 1
3  sample_size = df_contraceptives.shape[0]
4  beta_0 = np.random.normal(beta_0_mu, beta_0_sigma, sample_size)
5  beta_1 = 4
6  beta_2 = -3
7  beta_3 = -2
8
9  df_contraceptives_simulated = df_contraceptives.copy()
10 df_contraceptives_simulated['contraceptive_use_sim'] = expit(beta_0 +
11                                     beta_1*df_contraceptives_simulated['urban'] +
12                                     beta_2*df_contraceptives_simulated['living.children'] +
13                                     beta_3*df_contraceptives_simulated['age_mean'])
14 df_contraceptives_simulated['contraceptive_use_sim'] = (df_contraceptives_simulated['contraceptive_use_sim'] > 0.5).astype(int)
15 df_contraceptives_simulated.head()
```

bayes\_logreg\_simulated.py hosted with ❤ by GitHub

[view raw](#)

district urban living.children age mean contraceptive use contraceptive use sim



Get started

Open in app



29	0	2	-8.5599	1	1
5	0	3	-4.5599	1	1
34	1	4	8.4400	0	0

## Does contraceptives usage vary by district?

In order to experiment with the Bayesian approach a bit more, we will now specify a varying-intercept logistic regression model, where the intercept varies by district, and we will fit it to the simulated contraceptives data.

We set informative prior distributions on  $\beta_0$ ,  $\beta_1$ ,  $\beta_2$  and  $\beta_3$ , which are explicitly different from the ones we used to generate the simulated contraceptives data. Our goal is to estimate the true parameter values, which were used to simulate the response variable  $Y$ .

### Prior Distribution:

$\beta_{0j} \sim N(\mu_0, \sigma_0^2)$ , with  $\mu_0 \sim N(0, 10000)$  and  $\frac{1}{\sigma_0^2} \sim \text{Gamma}(0.1, 0.1)$ .

$\beta_1 \sim N(0, 10000)$ ,  $\beta_2 \sim N(0, 10000)$ ,  $\beta_3 \sim N(0, 10000)$

### Generative Model for data:

$Y_{ij} \sim \text{Bernoulli}(p_{ij})$

$\text{logit } p_{ij} = \beta_{0j} + \beta_1 \times \text{urban} + \beta_2 \times \text{living.children} + \beta_3 \times \text{age-mean}$ ,

where  $Y_{ij}$  is 1 if woman  $i$  in district  $j$  uses contraceptives, and 0 otherwise, and where  $i \in \{1, \dots, N\}$  and  $j \in \{1, \dots, J\}$ .  $N$  is the number of observations in the data, and  $J$  is the number of districts. The above notation assumes  $N(\mu, \sigma^2)$  is a Normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

When we have data, priors and a generative model, we can apply Bayes theorem to compute the **posterior probability distribution** of the model parameters conditionally upon the predictors (*district*, *urban*, *living.children*, *age-mean*) and response ( $Y$ ).

*Markov chain Monte Carlo (MCMC)* is a popular class of algorithms used to find the posterior distribution of the model parameters. Running such algorithms in Python is

[Get started](#)[Open in app](#)

chains to guide the process towards the simulated data we have at hand. You can see how the MCMC works in the code below.

```
1  import pyjags
2
3  # specs of varying-intercept model in PyJags syntax
4  logreg_code = '''
5  model {
6      for (i in 1:N){
7          Y[i] ~ dbern(p[i])
8          p[i] = ilogit(beta_0[district[i]] + beta_1*urban[i] + beta_2*children[i] + beta_
9      }
10     for (j in 1:J){
11         beta_0[j] ~ dnorm(beta_0_mu, beta_0_sigma)
12     }
13     beta_0_mu ~ dnorm(0, 1/10000)
14     beta_0_sigma ~ dgamma(0.1, 1/pow(0.1,2))
15     beta_1 ~ dnorm(0, 1/10000)
16     beta_2 ~ dnorm(0, 1/10000)
17     beta_3 ~ dnorm(0, 1/10000)
18 }
19 '''
20
21 # number of Markov chains
22 nb_chains = 3
23
24 # builder for PyJags model
25 def get_pyjags_model(Y):
26     logreg_data = dict(district=df_contraceptives_simulated['district'],
27                       urban=df_contraceptives_simulated['urban'],
28                       children=df_contraceptives_simulated['living.children'],
29                       age=df_contraceptives_simulated['age_mean'],
30                       Y=Y,
31                       N=sample_size,
32                       J=nb_districts)
33     logreg_init = dict( beta_0_mu=beta_0_mu,
34                       beta_0_sigma=beta_0_sigma,
35                       beta_1=beta_1,
36                       beta_2=beta_2,
37                       beta_3=beta_3
38     )
```



Get started

Open in app

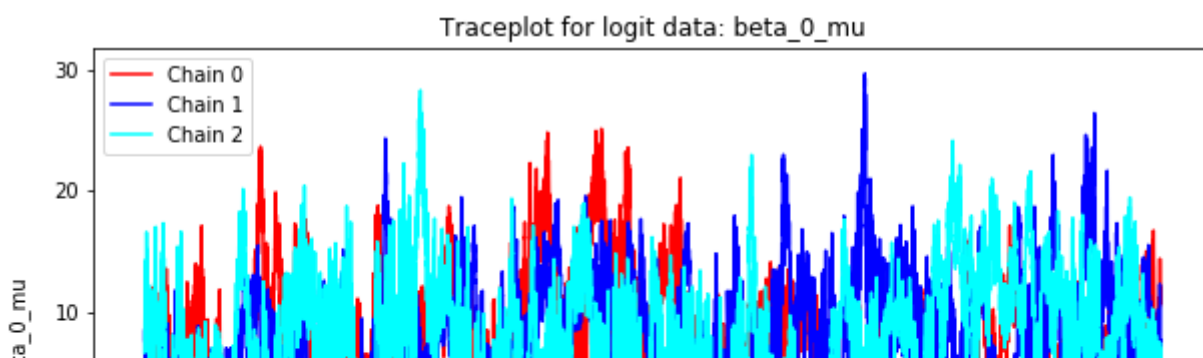


```

42 # varying-intercept PyJags model
43 logreg_model = get_pyjags_model(Y=df_contraceptives_simulated['contraceptive_use_sim'])
44
45 # sampler function
46 def get_pyjags_samples(logreg_model=logreg_model, n_burnin=1000, n_samples=10000):
47     logreg_burnin = logreg_model.sample(n_burnin)
48     logreg_samples = logreg_model.sample(n_samples)
49     return logreg_burnin, logreg_samples
50
51 # sampling from posterior distribution
52 logreg_burnin, logreg_samples = get_pyjags_samples()
53
54 # list of parameters
55 logreg_parameters = ['beta_0_mu', 'beta_0_sigma', 'beta_1', 'beta_2', 'beta_3']
56
57 # function for plotting the posterior distributions trace
58 def trace_plot(logreg_samples=logreg_samples, n_samples=10000):
59     colors = ['red', 'blue', 'cyan']
60     fig, ax = plt.subplots(len(logreg_parameters), 1, figsize=(10,30))
61     for i, param in enumerate(logreg_parameters):
62         for j in range(nb_chains):
63             ax[i].plot(range(n_samples),
64                         logreg_samples[param][0,:,j], color=colors[j], label='Chain {}'.for
65                         )
66             ax[i].set_xlabel('iteration')
67             ax[i].set_ylabel(param)
68             ax[i].set_title('Traceplot for logit data: {}'.format(param))
69             ax[i].legend()
70
71 trace_plot()

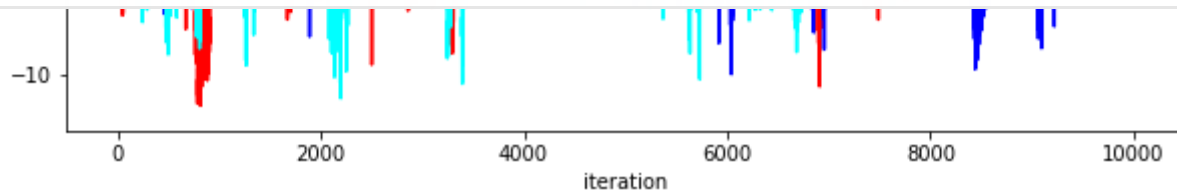
```

bayes\_logreg\_pyjags1.py hosted with ❤ by GitHub

[view raw](#)

Get started

Open in app



A good practice is to plot the trace plots of the MCMC sampler for the parameters. The trace plot above suggests that the sampler did converge. The average of each of the 3 Markov chains looks roughly the same. Although there is a little wandering within the chain, there is no evidence of divergent chains.

Now, we can plot the histograms of the posterior distributions for the parameters. We only show the plot for  $\beta_0$  and for district10 only, as illustrative example.

```

1  # function to extract posterior distributions
2  def get_pyjags_dataframe(logreg_samples=logreg_samples,
3                          logreg_parameters=logreg_parameters,
4                          nb_chains=nb_chains):
5      df_samples = []
6      for chain in range(nb_chains):
7          df_chain = []
8          for param in logreg_parameters:
9              samples = logreg_samples[param]
10             chain_data = samples[:, :, chain]
11             df = pd.DataFrame(chain_data.T)
12             if df.shape[1]==1:
13                 df = df.rename({0:param}, axis=1)
14             else:
15                 df = df.add_prefix(param)
16             df_chain.append(df)
17
18             df_samples.append(pd.concat(df_chain, axis=1))
19
20     return df_samples
21
22 # get posterior distributions of beta0 for all districts
23 beta_0_samples = get_pyjags_dataframe(logreg_samples, ['beta_0'], nb_chains)
24
25 # plot posterior distributions of beta0 for selected districts
26 selected_districts = [10,20,30,40,50,60]
```

Get started

Open in app

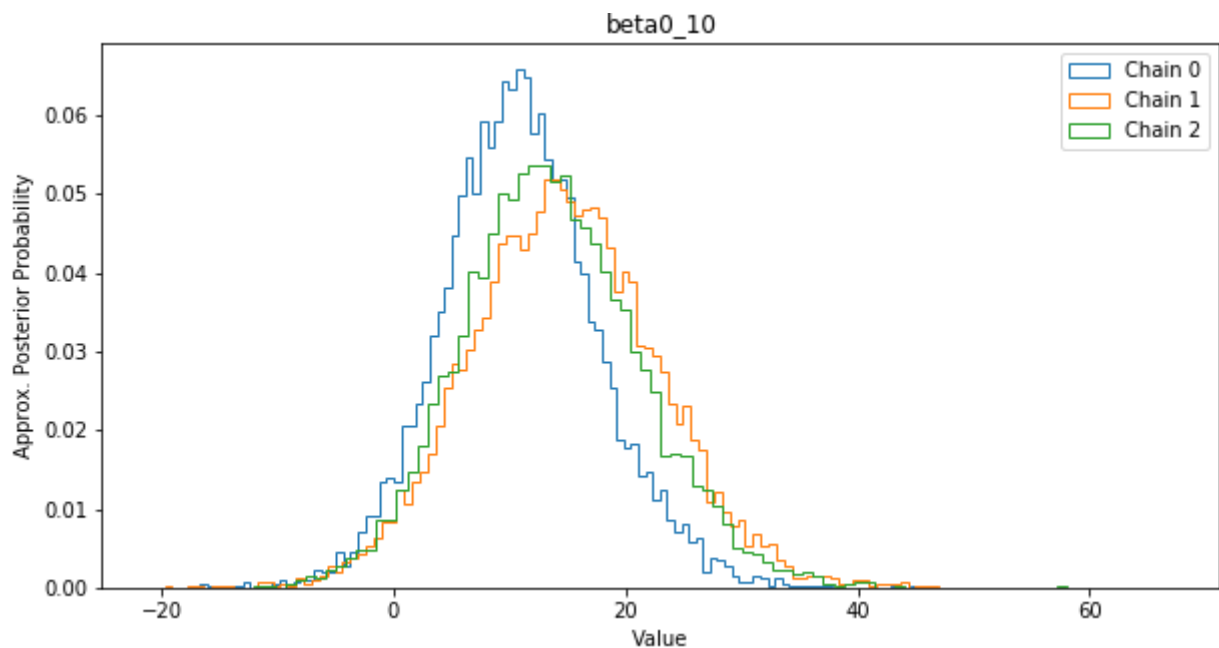


```

30     col_idx = int(col[-2:])+1
31     if col_idx in selected_districts:
32         for chain in range(nb_chains):
33             ax[i].hist(beta_0_samples[chain][col],
34                       bins=100, histtype='step', density=True, label="Chain {}".format(
35             ax[i].set_xlabel("Value")
36             ax[i].set_ylabel("Approx. Posterior Probability")
37             ax[i].set_title('beta0_{}'.format(col_idx))
38             ax[i].legend()
39     i = i + 1

```

bayes\_logreg\_pyjags\_posteriors1.py hosted with ❤ by GitHub

[view raw](#)

We recall that the true distribution for  $\beta_0$  that was used to generate simulated data was as follows.

$$\beta_0 \sim N(2, 1)$$

As you can see in the plot above, the true  $\beta_0$  parameter for district 10 is contained within the posterior distributions from our model. In fact, this was also the case for all remaining parameters, not shown here. This finding suggests that the Bayesian

Get started

Open in app

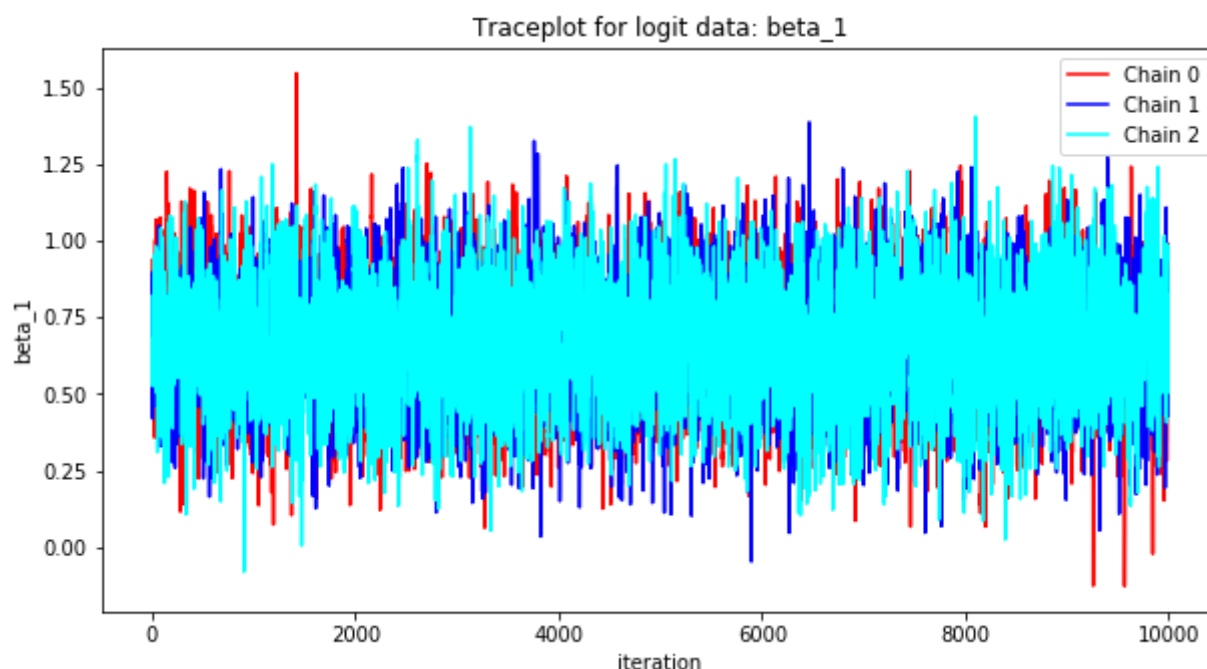


## Women belonging to which district are most likely to use contraceptives?

Below, we run the MCMC sampler once more, this time using training data. We check the convergence by examining the trace plots, as we did with the simulated data. It appears that the samplers also converged here. None of the chains appears to be divergent, because they look uncorrelated, independently random sampled.

```
1 # create PyJAGS model
2 logreg_model_train = get_pyjags_model(Y=df_contraceptives_simulated['contraceptive_use'])
3
4 # get samples
5 logreg_train_burnin, logreg_train_samples = get_pyjags_samples(logreg_model=logreg_model_
6
7 # plot trace
8 trace_plot(logreg_samples=logreg_train_samples)
```

bayes\_logreg\_pyjags2.py hosted with ❤ by GitHub

[view raw](#)

The posterior distribution of model parameters  $\beta_0$  can now be plotted for all districts as follows.

Get started

Open in app

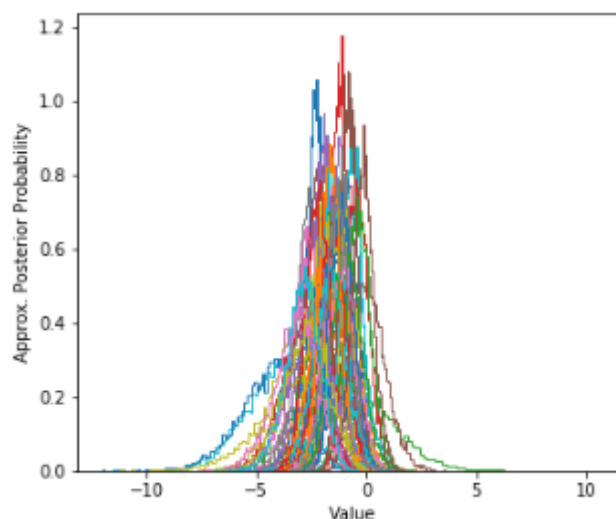


```

4                                     nb_chains=nb_chains)
5
6 # plot posterior distributions of beta0 for all districts
7 beta_0_avg = []
8 beta_0_std = []
9 fig, ax = plt.subplots(1, nb_chains, figsize=(20, 5))
10 for chain in range(nb_chains):
11     col_avg = []
12     col_std = []
13     for col in beta_0_samples_train[0].columns:
14         col_avg.append(np.mean(beta_0_samples_train[chain][col]))
15         col_std.append(np.std(beta_0_samples_train[chain][col]))
16         ax[chain].hist(beta_0_samples_train[chain][col],
17                        bins=100, histtype='step', density=True, label=col)
18         ax[chain].set_xlabel("Value")
19         ax[chain].set_ylabel("Approx. Posterior Probability")
20         ax[chain].set_title('Chain {}'.format(chain))
21     beta_0_avg.append(col_avg)
22     beta_0_std.append(col_std)

```

bayes\_logreg\_pyjags\_posteriors2.py hosted with ❤ by GitHub

[view raw](#)

We know that positive values of  $\beta_0$  are associated with increased probability that women belonging to the corresponding districts are most likely to use contraceptives. Looking at the plot above, only a few districts have positive  $\beta_0$  values.

Get started

Open in app



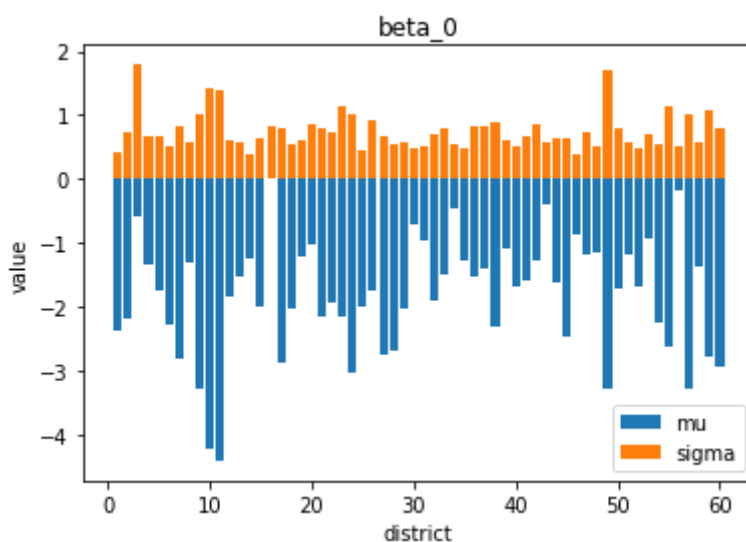
group.

```

1  beta_0_avg_arr = np.mean(np.array(beta_0_avg), axis=0)
2  beta_0_std_arr = np.mean(np.array(beta_0_std), axis=0)
3  print('Women belonging to district {} are most likely to use contraceptives.'.format(np.
4  print('Women belonging to district {} are most likely NOT to use contraceptives.'.format
5
6  plt.bar(range(1,61),beta_0_avg_arr, label='mu')
7  plt.bar(range(1,61),beta_0_std_arr, label='sigma')
8  plt.xlabel('district')
9  plt.ylabel('value')
10 plt.title('beta_0')
11 plt.legend();

```

bayes\_logreg\_pyjags\_evaluation1.py hosted with ❤ by GitHub

[view raw](#)

The posterior distributions allow us to compute the mean and standard deviation of  $\beta_0$  values per district as illustrated in the plot above. Based on these results, we can conclude that the  $\beta_0$  posteriors are not distributed uniformly across districts. This evidence is in support of the varying-intercept model and lead to the following findings:

*Women belonging to district 16 are most likely to use contraceptives.*

*Women belonging to district 11 are most likely NOT to use contraceptives.*

## Conclusion



[Get started](#)[Open in app](#)

When combined with prior beliefs, we were able to quantify uncertainty around point estimates of contraceptives usage per district. In this article, we also offered few take-out on PyJAGS, an easy to use Python library for Bayesian inference.

This area opens a wide door for future work, especially because Bayesian statistical analysis is at the core of several technologies ranging from medical diagnosis to election forecasting. The *Signal and the Noise* 2012's book by [Nate Silver](#) is an example of master piece in the art of using probability and statistics as applied to real-world circumstances.

Thanks to Anne Bonner.

## Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

[Get this newsletter](#)[Machine Learning](#)[Data Science](#)[Artificial Intelligence](#)[Statistics](#)[Towards Data Science](#)[About](#) [Write](#) [Help](#) [Legal](#)

Get the Medium app

