

Supplementary Materials for Introducing Variational Inference In Undergraduate Statistics and Data Science Curriculum

Vojtech Kejzlar

Department of Mathematics and Statistics, Skidmore College
and

Jingchen Hu

Department of Mathematics and Statistics, Vassar College

January 2, 2023

Abstract

The supplementary materials include: 1) Details of the class activity on probabilistic model for count data with variational inference, introduced in Section 3 in the main text; 2) The manual of the R `shiny` app we have developed for the module, mentioned in Section 3 in the main text; and 3) Details of the guided R lab of the LDA application to a sample of the Associated Press newspaper articles with variational inference, presented in Section 4 in the main text.

Keywords: Active learning, Bayesian Statistics, Computing, MCMC, Probabilistic Models, Undergraduate Curriculum

1 Class Activity: Probabilistic Model for Count Data with Variational Inference

The goal of this activity is to illustrate variational inference on a simple example of Gamma-Poisson conjugate model, which is a popular model for count data.

1.1 A Motivating Example

Our task is to estimate the average number of active users of a popular massively multiplier online role-playing game (mmorpg) playing between the peak evening hours 7 pm and 10 pm. This information can help the game developers in allocating server resources and optimizing user experience. To make this estimate, we will consider the following counts (in thousands) of active players collected during the peak evening hours over a two-week period past month.

	Sun	Mon	Tue	Wed	Thu	Fri	Sat
Week 1	50	47	46	52	49	55	53
Week 2	48	45	51	50	53	46	47

1.2 Overview of the Gamma-Poisson Model

Sampling density:

Suppose that $\mathbf{y} = (y_1, \dots, y_n)$ represent the observed counts in n time intervals where the counts are independent, then each y_i follows a Poisson distribution with rate $\theta > 0$. Namely,

$$y_i \mid \theta \sim \text{Poisson}(\theta)$$

- $\mathbb{E}(y_i \mid \theta) = \theta$
- $\text{Var}(y_i \mid \theta) = \theta$

Prior distribution:

$$\theta \sim \text{Gamma}(\alpha, \beta)$$

- $\alpha > 0$ is the shape parameter
- $\beta > 0$ is the rate parameter
- $\mathbb{E}(\theta) = \frac{\alpha}{\beta}$
- $\mathbb{V}ar(\theta) = \frac{\alpha}{\beta^2}$

Posterior distribution:

$$\theta \mid y_1, \dots, y_n \sim \text{Gamma}(\alpha + \sum_{i=1}^n y_i, \beta + n)$$

1.3 Exact Inference with the Gamma-Poisson Model

We will start by selecting a prior distribution for the unknown average number of active users. Suppose that we elicit an expert's advice on the matter, and they tell us that a similar mmorpg has typically about 50,000 users during peak hours. However, they are not too sure about that, so the interval between 45,000 and 55,000 users should have a reasonably high probability. Suppose that we elicit an expert's advice on the matter, and they tell us that a similar mmorpg has typically about 50,000 users during peak hours. However, they are not too sure about that, so the interval between 45,000 and 55,000 users should have a reasonably high probability. This reasoning leads to a $\text{Gamma}(100, 2)$ as a reasonable prior for the average number of active users.

Task 1: Explain the reasoning behind using $\text{Gamma}(100, 2)$ as the prior distribution.

Task 2: Use the information above to find the exact posterior distribution for the average number of active users.

Task 3: What are the mean and standard deviation of the posterior distribution that you just obtained? What is your recommendation about the typical number of active users playing the mmorpg between the peak evening hours 7pm and 10pm?

1.4 Variational Inference with the Gamma-Poisson Model

Variational inference approximates the (unknown) posterior distribution of a parameter by a simple family of distributions. In this case, we will try to approximate the posterior distribution of the mmorpg's average number of active users between the peak hours θ by a log-normal distribution with mean μ and standard deviation σ . Log-normal distribution is a continuous probability distribution of a random variable whose logarithm is normally distributed. It also happens to be a popular variational family for non-negative parameters as it is amenable to autodifferentiation. Since we know exactly how the posterior distribution for Gamma-Poisson model looks like, we will be able to check the fidelity of the variational approximation. Use the accompanying applet titled *Variational Inference with Gamma-Poisson Model for count data* to complete the following task.

Task 4: Use the sliders in the applet to manually find the member of a log-normal variational family that well approximates the posterior distribution of θ . What is your strategy?

Task 5: Compare your approximation with a neighbor. Whose approximation is closer to the exact posterior distribution of θ ? How are you deciding?

Task 6: Check the *Fit a variational approximation* box in the applet to find the variational approximation using the gradient ascent algorithm. How close was the variational approximation that you found manually to the one found here?

Variational Inference with Gamma-Poisson Model for count data

Variational approximation using log-normal variational family:

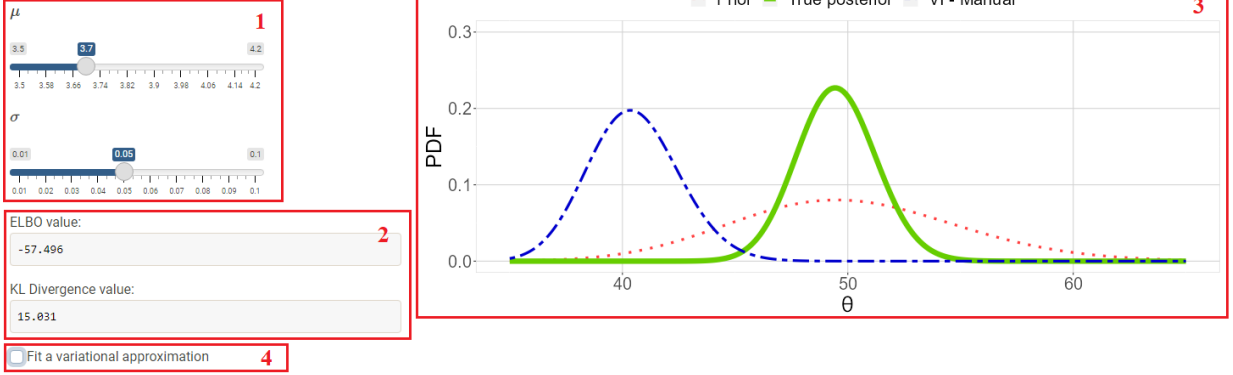


Figure 1: The applet is based on the class activity presented in Section 1 of the supplementary materials. The applet visual before checking the “Fit a variational approximation” checkbox is displayed.

2 Manual of the R shiny app

This document describes the elements of R Shiny applet that accompanies the “Probabilistic Model for Count Data with Variational Inference” class activity. Note that the numbering in Section 2.1 and Section 2.2 corresponds to the numbered boxes in Figure 1 and Figure 2.

2.1 Manual search for variational approximation

1. Sliders to control the mean μ and the standard deviation σ of log-normal variational family.
2. The ELBO and KL divergence values for variational approximation based on the mean and standard deviations selected in box 1.
3. A plot that displays the true Gamma(792, 100) posterior distribution, the Gamma(100, 2) prior distribution, and the variational approximation based on the selection in box 1.
4. A checkbox to display the results of ELBO maximization via gradient ascent algorithm. The resulting variational approximation is plotted in box 3.

Variational Inference with Gamma-Poisson Model for count data

Variational approximation using log-normal variational family:

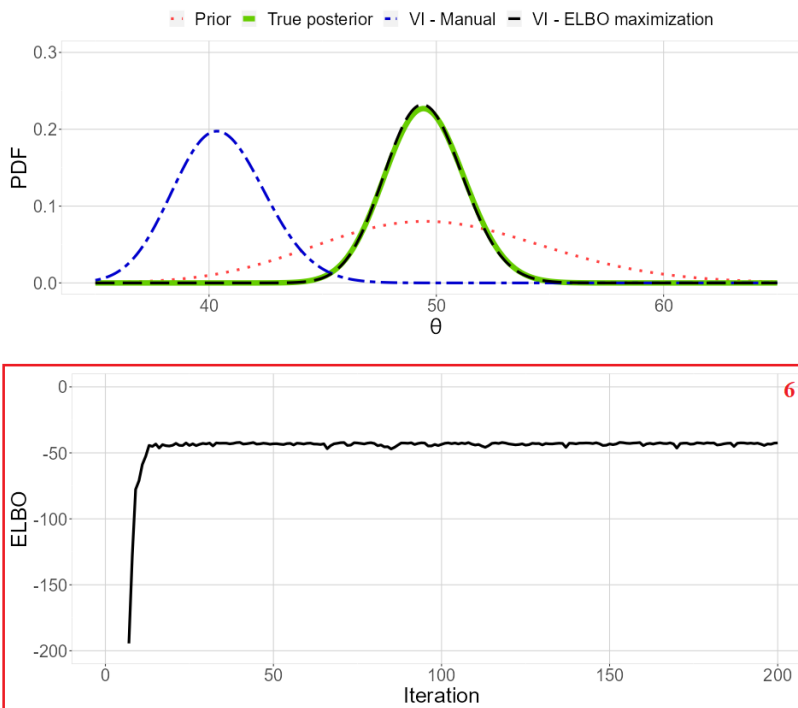
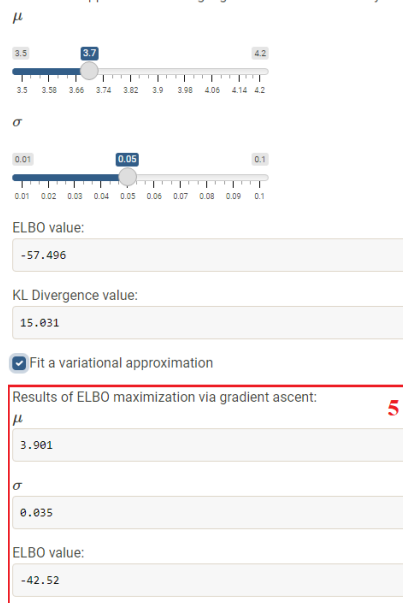


Figure 2: The applet visual after checking the “Fit a variational approximation“ checkbox is displayed.

2.2 Variational approximation based on ELBO maximization

5. The resulting mean μ , standard deviation σ , and ELBO values of variational approximation based on ELBO maximization.
6. A plot depicting ELBO values for each iteration of the gradient ascent algorithm.

3 Lab: Document Clustering

The goal of this lab is to gain a practical experience with variational inference on a realistic use case based on Latent Dirichlet Allocation (LDA) and implement the model in R applied to a dataset of documents. To do so, we consider a collection of 2246 Associated Press newspaper articles to be clustered using the LDA model. The dataset is part of the `topicmodels` R package. You can load the dataset `AssociatedPress` with the following R

command.

```
data("AssociatedPress", package = "topicmodels")
```

3.1 Overview of the LDA model and Stan script

The LDA is a mixed-membership clustering model, commonly used for document clustering. LDA models each document to have a mixture of topics, where each word in the document is drawn from a topic based on the mixing proportions. Specifically, the LDA model assumes K topics for M documents made up of words drawn from V distinct words. For document m , a topic distribution $\boldsymbol{\theta}_m$ is drawn over K topics from a Dirichlet distribution,

$$\boldsymbol{\theta}_m \sim \text{Dirichlet}(\boldsymbol{\alpha}), \quad (1)$$

where $\sum_{k=1}^K \theta_{m,k} = 1$ ($0 \leq \theta_{m,k} \leq 1$) and $\boldsymbol{\alpha}$ is the prior a vector of length K with positive values.

Each of the N_m words $\{w_{m,1}, \dots, w_{m,N_m}\}$ in document m is then generated independently conditional on $\boldsymbol{\theta}_m$. To do so, first, the topic $z_{m,n}$ for word $w_{m,n}$ in document m is drawn from

$$z_{m,n} \sim \text{categorical}(\boldsymbol{\theta}_m), \quad (2)$$

where $\boldsymbol{\theta}_m$ is the document-specific topic-distribution defined in Equation (1).

Next, the word $w_{m,n}$ in document m is drawn from

$$w_{m,n} \sim \text{categorical}(\boldsymbol{\phi}_{z_{m,n}}), \quad (3)$$

which is the word distribution for topic $z_{m,n}$. Note that $z[m, n]$ in Equation (3) refers to $z_{m,n}$.

Lastly, a Dirichlet prior is given to distributions $\boldsymbol{\phi}_k$ over words for topic k as

$$\boldsymbol{\phi}_k \sim \text{Dirichlet}(\boldsymbol{\beta}), \quad (4)$$

where $\boldsymbol{\beta}$ is the prior a vector of length V (i.e., the total number of words) with positive values. Figure 3 shows a graphical model representation of LDA.

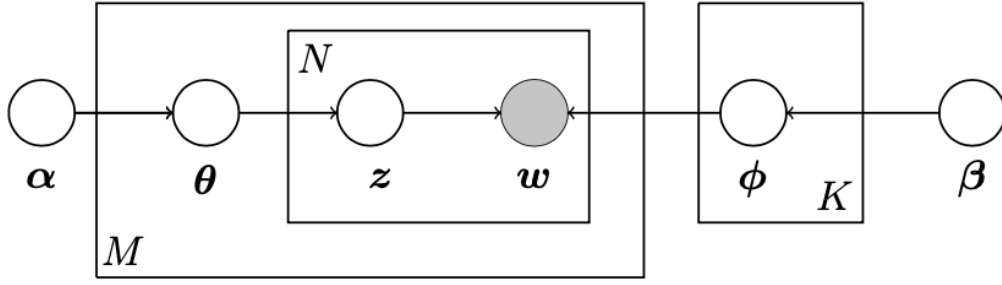


Figure 3: Graphical model representation of LDA. The outer box represents the documents. On the left, the inner box represents the topics and words within each document. On the right, the box represents the topics.

Below, we include the **Stan** script for the LDA model provided by Stan Development Team (2022). Note that **Stan** supports the calculation of marginal distributions over the continuous parameters by summing out the discrete parameters in mixture models (Stan Development Team 2022). This process corresponds to the **gamma** parameter in the **Stan** script below. We refer interested readers to Stan Development Team (2022) for further details.

```
data {
  int<lower=2> K;           // number of topics
  int<lower=2> V;           // number of words
  int<lower=1> M;           // number of docs
  int<lower=1> N;           // total word instances
  int<lower=1, upper=V> w[N]; // word n
  int<lower=1, upper=M> doc[N]; // doc ID for word n
  vector<lower=0>[K] alpha; // topic prior vector of length K
  vector<lower=0>[V] beta;  // word prior vector of length V
}

parameters {
  simplex[K] theta[M]; // topic distribution for doc m
}
```



```

    simplex[V] phi[K];      // word distribution for topic k
  }

model {
  for (m in 1:M)
    theta[m] ~ dirichlet(alpha);
  for (k in 1:K)
    phi[k] ~ dirichlet(beta);
  for (n in 1:N) {
    real gamma[K];
    for (k in 1:K)
      gamma[k] = log(theta[doc[n], k]) + log(phi[k, w[n]]);
    target += log_sum_exp(gamma); // likelihood;
  }
}

```

3.2 Variational inference with the LDA model

For demonstration purposes, we shall start with a two-topic LDA model (i.e., $K = 2$). Before that, we recommend removing the words from **AssociatedPress** datasets that are rare using the function `removeSparseTerms()` from the `tm` package. These words have a minimal effect on the LDA parameter estimation. Nevertheless, they increase the computational cost of variational inference and therefore should be removed using the following R command.

```
dtm <- removeSparseTerms(AssociatedPress, 0.95)
```

We are now ready to fit the LDA model using variational inference capabilities of the `cmdstanr` package. The following code achieves the goal:

```

LDA_model_cmd <- cmdstan_model(stan_file = "LDA.stan")

N_TOPICS <- 2

data <- list(K = N_TOPICS,
            V = dim(dtm)[2],
            M = dim(dtm)[1],
            N = sum(dtm$v),
            w = rep(dtm$j, dtm$v),
            doc = rep(dtm$i, dtm$v),
            #according to Griffiths and Steyvers(2004)
            alpha = rep(50/N_TOPICS, N_TOPICS),
            beta = rep(1, dim(dtm)[2])
)

vi_fit <- LDA_model_cmd$variational(data = data,
                                   seed = 1,
                                   output_samples = 1000,
                                   eval_elbo = 1,
                                   grad_samples = 10,
                                   elbo_samples = 10,
                                   algorithm = "meanfield",
                                   output_dir = NULL,
                                   iter = 1000,
                                   adapt_iter = 20,
                                   save_latent_dynamics=TRUE,
                                   tol_rel_obj = 10^-4)

```

The “LDA.stan” file contains the **Stan** script for the LDA model provided in Section 3.1. We recommend the usage of the R help to get familiar with the `variational()` method of the `cmdstan_model()` function. The variable `vi_fit` contains the results of variational approximation of the LDA parameters. For example, one can obtain the word distributions

for the each of the topics with `vi_fit$summary("phi")`.

Finally, to access the ELBO values, use the following:

```
vi_diag <- utils::read.csv(vi_fit$latent_dynamics_files()[1],  
                           comment.char = "#")  
ELBO <- data.frame(Iteration = vi_diag[,1], ELBO = vi_diag[,3])
```

Task 1: Use a graphical display to show the 10 most common words for each of the two-topics and their probabilities.

Task 2: Use the function `wordcloud()` from the `wordcloud` package and display the most common words for each of the topics as a word cloud. What kinds of articles do these topics represent?

Task 3: Fit a three-topic LDA model, display the most common words for each of the topics. How do the results differ from the two-topic LDA?

Task 4 (Advanced): Use the three-topic LDA model and display the topic prevalence among the 2246 Associated Press articles. That is, show what proportions of articles fall under each topic.

All necessary R code for fitting the LDA model to the Associated Press sample, including the graphical displays shown in the main text, is included in a separate R script file called `LDA.LAB.R` available as a part of the supplementary materials. We also include a printout of the R script below for interested readers.

```
library(cmdstanr)  
  
# Checking integrity of installation of cmdstanr  
check_cmdstan_toolchain()  
install_cmdstan(cores = 2)  
cmdstan_path()  
cmdstan_version()
```

```

# Auxiliary packages
library(tm)
library(tidyverse)
library(tidytext)
library(topicmodels)

## Get data
data("AssociatedPress", package = "topicmodels")

## Removing rare words from the vocabulary
dtm <- removeSparseTerms(AssociatedPress, 0.95)
dim(dtm)

## Input for stan model
N_TOPICS <- 2

data <- list(K = N_TOPICS,
             V = dim(dtm)[2],
             M = dim(dtm)[1],
             N = sum(dtm$v),
             w = rep(dtm$j, dtm$v),
             doc = rep(dtm$i, dtm$v),
             #according to Griffiths and Steyvers(2004)
             alpha = rep(50/N_TOPICS, N_TOPICS),
             beta = rep(1, dim(dtm)[2])
)

### VB fit
LDA_model_cmd <- cmdstan_model(stan_file = "LDA.stan")
LDA_model_cmd$print()

```

```

vb_fit <- LDA_model_cmd$variational(data = data,
                                     seed = 1,
                                     output_samples = 1000,
                                     eval_elbo = 1,
                                     grad_samples = 10,
                                     elbo_samples = 10,
                                     algorithm = "meanfield",
                                     output_dir = NULL,
                                     iter = 1000,
                                     adapt_iter = 20,
                                     save_latent_dynamics=TRUE,
                                     tol_rel_obj = 10^-4)

# Plotting ELBO
vb_diag <- utils::read.csv(vb_fit$latent_dynamics_files()[1],
                           comment.char = "#")
ELBO <- data.frame(Iteration = vb_diag[,1],
                  ELBO = vb_diag[,3])

ggplot(data = ELBO, aes(x = Iteration, y = ELBO)) + geom_line(lwd=1.5) +
  theme(text = element_text(size = 20),
        panel.background = element_rect(fill = "transparent",
                                          color = "lightgrey"),
        panel.grid.major = element_line(colour = "lightgrey")) +
  xlim(0,110)

## Accessing parameters
vb_fit$summary("theta") # dim: M-by-K
vb_fit$summary("phi") # dim: K-by-V

```

```

## Word distribution per topic
V <- dim(dtm)[2]
odd_rows <- rep(c(1,0), times = V)
Topic1 <- vb_fit$summary("phi")[odd_rows == 1,]
Topic2 <- vb_fit$summary("phi")[odd_rows == 0,]

word_probs <- data.frame(Topic = c(rep("Topic 1", V),
                                   rep("Topic 2", V)),
                        Word = rep(dtm$dimnames$Terms, N_TOPICS),
                        Probability = c(Topic1$mean, Topic2$mean))

# Selecting top 10 words per topic
top_words <- word_probs %>% group_by(Topic) %>% top_n(10) %>%
  ungroup() %>% arrange(Topic, -Probability)

top_words %>%
  mutate(Word = reorder_within(Word, Probability, Topic)) %>%
  ggplot(aes(Probability, Word, fill = factor(Topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ Topic, scales = "free") +
  scale_y_reordered() + theme(text = element_text(size = 15)) + xlim(0,0.025) +
  xlab("Word distributions ( \u03d5 )")

# Word Cloud display
#install.packages("wordcloud")
library(wordcloud)

top_words <- word_probs %>% group_by(Topic) %>% top_n(20) %>%
  ungroup() %>% arrange(Topic, -Probability)

```

```

mycolors <- brewer.pal(8, "Dark2")
wordcloud(top_words %>% filter(Topic == "Topic 1") %>% .$Word ,
          top_words %>% filter(Topic == "Topic 1") %>% .$Probability,
          random.order = FALSE,
          color = mycolors)

mycolors <- brewer.pal(8, "Dark2")
wordcloud(top_words %>% filter(Topic == "Topic 2") %>% .$Word ,
          top_words %>% filter(Topic == "Topic 2") %>% .$Probability,
          random.order = FALSE,
          color = mycolors)

```

References

Stan Development Team (2022), *Stan Modeling Language User's Guide and Reference Manual*, Version 2.31.
URL: <http://mc-stan.org/>