

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

---

# Introduction to Bayesian Modeling

---

## Contents

<b>1 Overview</b>	<b>2</b>
1.1 Goal . . . . .	2
1.2 Date and Time . . . . .	2
1.3 Target Audience . . . . .	2
1.4 Prerequisites . . . . .	2
1.5 Textbook . . . . .	2
1.6 Philosophy . . . . .	2
1.7 Format . . . . .	2
<b>2 Topics</b>	<b>3</b>
2.1 Introduction to Bayes . . . . .	3
2.2 Methods . . . . .	3
2.3 More complicated models . . . . .	4
<b>A Resources which may be appropriate for mini-reading group or student mini-presentations</b>	<b>5</b>
A.1 Intro . . . . .	5
A.2 Model Checking . . . . .	5
A.3 Missing data and imputation . . . . .	5
<b>B Mini projects</b>	<b>5</b>
B.1 Batting average dataset . . . . .	5
<b>C Notes to self: some details for when slides are constructed</b>	<b>6</b>
C.1 Why Bayes . . . . .	6
C.2 Introduction to inference . . . . .	6
C.3 Hierarchical models . . . . .	6

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# 1 Overview

## 1.1 Goal

The goal of this workshop is to introduce students to the concepts and practice of Bayesian modeling. We will begin by motivating Bayesian approaches. Next, we will introduce and apply models with conjugate priors, such as Bayesian normal, Bayesian binomial, and Bayesian linear regression. We will then introduce the two primary techniques for approximate Bayesian inference, namely Markov Chain Monte Carlo (MCMC) and variational inference. Using these techniques, and in some cases clever trickery, we will then tackle models for which there are not conjugate priors, such as Bayesian logistic regression, Bayesian multiclass regression, Bayesian mixture models, and Bayesian hidden Markov models. Finally, we will very briefly discuss Bayesian deep learning. For applications, we will use Python; namely, a combination of pymc3, scikit-learn, and code we write ourselves.

## 1.2 Date and Time

The workshop will be held via Zoom, June 7-11, from 2pm-5pm EST.

## 1.3 Target Audience

We expect that the typical student will be a graduate student, faculty member, staff member, or researcher in a quantitative field (such as computer science, statistics, engineering, or biology), who would like to learn more about Bayesian modeling.

## 1.4 Prerequisites

Prerequisites include calculus, some familiarity with introductory linear algebra (matrix multiplications, determinants, and traces), and some familiarity with introductory probability (e.g., we will assume prior fluency with concepts such as expectation, conditional probability, and commonly used distributions, such as Gaussian and Poisson). The class will use Python as a common language. The workshop will employ student-centered components; for maximum benefits, we strongly encourage setting aside 1-2 hours per day outside of the workshop to work on material.

## 1.5 Textbook

We will use [1], available online at <http://www.stat.columbia.edu/~gelman/book/>.

## 1.6 Philosophy

*Learning in order to create* is both more fun and more effective than *learning for some extrinsic purpose*. Hence, the workshop is structured so as to (a) be student-centered and (b) allow self-determination and autonomy in how students engage with the material.

## 1.7 Format

About half of the workshop will involve lectures via slides.

About half of the workshop will be interactive, including:

- Student “lightning chat” (10 minute) presentations. Something like one per student per workshop, depending on the number of students. To allow for student-centered direction and autonomy, students may choose any of the following:
  - Presentation of Python implementations of models from [2], [1], or the workshop.
  - Presentation of an exercise from [1] or [2].
  - Reviewing a demo with us from [https://github.com/avehtari/BDA\\_py\\_demos](https://github.com/avehtari/BDA_py_demos).
  - Presentation of a reading section, blog, etc. of interest.
  - Presentation of a mathematical derivation of something relevant to the course.

- Presentation of how a concept relates to something from their research area.
- Real-time python applications lab – Google Collab exercises ? Python (rather than R) implementations of [2] and [1].
- Mini reading group discussions – They might not have time to read a whole paper, but we could discuss sections of a text at the very beginning, or perhaps sections of a relevant paper.

## 2 Topics

Below are topics we plan to cover in the course:

### 2.1 Introduction to Bayes

We present everything in here using conjugate models with closed-form posteriors. The models are useful in and of themselves, as well as to build intuition for more complicated models.

Primary references here are [2] and [1].

- **Why Bayes?** – See Section 1.2 of [2]. [3] has some nice plots motivating why use Bayesian linear regression over standard linear regression. [4] has some nice plots illustrating the Bayesian approach and how it mitigates overfitting. I can provide a nice example with biometric profiling of human typing dynamics. [5] has a nice simple example of obtaining non-standard functionals from the posterior that can be of interest. [6] presents the case for Bayesian deep learning.
- **Belief functions, Bayes rule** – Sections 2.1, 2.2 of [2]. [4] briefly overviews of the Bayesian framework. For important mistakes in real life in medicine and law, see [7]. *Why most published research findings are false* [8] provides nice additional motivation in science. Could perhaps cover exchangeability here.
- **Introduction to inference** Chapter 11. 2 of [9] has a nice very brief introduction to inference. Sections 3.1, 3.2, 5, and 5 of [2] covers binomial, Poisson, normal, multivariate normal models. Introduce the exponential family formalism (see [10]; see also Section 5.2 of [9] ) for much greater breadth.
- **Bayesian linear regression** – Section 9 of [2]. I have notes on this. There are some nice slides here which also illustrate the use of kernels.<sup>1</sup> Introduce model selection here as well (Section 9.3 of [2] or Section 3.4 of [3]). Section 11.2.2 of [9] also has a nice, quick summary of Bayes factors for model comparison.
- **Bayesian workflow** – Lots of nice resources for Bayesian workflow. For example: [11] or [12]. Section 6 of [1] covers model checking, as does Section 11.2.3 of [9]. Some points to make re: model checking
  - *Samples from the posterior predictive should capture important properties of the observed dataset.* For a violation of this, see the normal model for Newcomb's speed of light measurements. (Compare Figures 6.2 and 3.1 of [1].)

We will want to find a way to get students to group up, probably based on domain expertise/interests, so that they can eventually work together on a project.

### 2.2 Methods

We introduce these methods, which can be used for models without closed-form posteriors. We practice applying them in the next section.

- **MCMC** - [Karin](#) will present, including an introduction to pymc3.
- **Variational inference** [13].

---

<sup>1</sup>Nice Bayesian linear regression slides: [https://www.cs.toronto.edu/~rgrosse/courses/csc411\\_f18/slides/lec19-slides.pdf](https://www.cs.toronto.edu/~rgrosse/courses/csc411_f18/slides/lec19-slides.pdf)

## 2.3 More complicated models

Here are some models which are still fairly standard, but lack conjugate priors, and so inference typically requires VI or MCMC. [Karin: Where here, or elsewhere, would you like to illustrate applications of MCMC?](#)

- **Hierarchical models** Section 11.4 of [9] has a nice brief overview. Hierarchical normal model (e.g. Gelman's 5 schools example). Hierarchical linear regression (Chapter 13 of [1], Secs 11.1-11.3 of [2]). Figure 11.1 and 11.3 (right) of [2] nicely shows the beneficial effect of sharing statistical strength in a hierarchical linear regression, as compared to many separate linear regressions.
- **Regression models for binary and multi-class data** Includes logistic regression, probit regression, binomial, multinomial, etc. Use this to cover additional inference techniques: auxilliary variable trick and Laplace variational inference. Show demo "beating scikit-learn with variational inference." See also pp. 390 of [2] for a useful warm-starting strategy. Could generalize to Bayesian GLMs. Could also cover or mention hierarchical extensions (i.e. Bayesian GLMM's).
- **Mixture models** I will give CAVI for Gaussian mixture models.
- **Time series models** Probably just hidden markov models<sup>2</sup>, although would be nice to also introduce state space models. Could mention embedding of GLM's or GLMM's within them. May give some overview to probabilistic graphical models here.
- **Bayesian Deep Learning** Bayes and neural networks. 20-30 min w/ guest presenter, Kyle Heuton, Ph.D. student, computer science.

## References

- [1] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC press, 2013.
- [2] Peter D Hoff. *A first course in Bayesian statistical methods*, volume 580. Springer, 2009.
- [3] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [4] Zoubin Ghahramani. Bayesian non-parametrics and the probabilistic approach to modelling. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 371(1984):20110553, 2013.
- [5] Leonhard Held and Chris C Holmes. Bayesian auxiliary variable models for binary and multinomial regression. *Bayesian analysis*, 1(1):145–168, 2006.
- [6] Andrew Gordon Wilson. The case for bayesian deep learning. *arXiv preprint arXiv:2001.10995*, 2020.
- [7] *The obscure maths theorem that governs the reliability of Covid testing*. Available at <https://www.theguardian.com/world/2021/apr/18/obscure-maths-bayes-theorem-reliability-covid-lateral-flow-tests-probability>.
- [8] John PA Ioannidis. Why most published research findings are false. *PLoS medicine*, 2(8):e124, 2005.
- [9] Anthony Christopher Davison. *Statistical models*, volume 11. Cambridge university press, 2003.
- [10] Michael Wojnowicz. *The exponential family*. Available (with permission).
- [11] Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C Margossian, Bob Carpenter, Yuling Yao, Lauren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. Bayesian workflow. *arXiv preprint arXiv:2011.01808*, 2020.
- [12] Jonah Gabry, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman. Visualization in bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(2):389–402, 2019.

<sup>2</sup>Nice reference for frequentist HMM's: <http://jwmi.github.io/ASM/5-HMMs.pdf>

- [13] Michael Wojnowicz. *Foundations of variational inference*. Available (with permission) at [https://github.com/mikewojnowicz/vi\\_foundations](https://github.com/mikewojnowicz/vi_foundations).
- [14] Ray Hill. Multiple sudden infant deaths—coincidence or beyond coincidence? *Paediatric and perinatal epidemiology*, 18(5):320–326, 2004.

## A Resources which may be appropriate for mini-reading group or student mini-presentations

### A.1 Intro

- Present *Why most published research findings are false*. [8].
- Present Bayesian analysis of multiple sudden infant deaths. [14].
- Textbook sections TBD.

### A.2 Model Checking

- Textbook sections TBD.

### A.3 Missing data and imputation

- Inference with a Bayesian normal model in the presence of missing data is given, in R code, on pp. 119 of [2]. Implement, and test, the code in Python.

## B Mini projects

These are small, open-ended problems to work on.

### B.1 Batting average dataset

The hierarchical normal model for (arcsine-transformed) batting average data on pp. 163 of [1] has some serious deficiencies, as exposed in Table 6.1 in the section on model checking.

Can you construct (and learn) a better model which makes predictions closer to the true final batting average?

Examples:

- Add an extra layer to the hierarchy, so that player  $p$ 's 1970 batting average inherits from player  $p$ 's overall batting average which in turn inherits from a population batting average. (Of course, I am speaking of the arcsine-transformed batting averages, so that we can use a hierarchical normal model.)
- Add an autoregressive component, because, as mentioned by Gelman, player batting averages *DO* change over time.

The text also does a poor job of checking the modeling assumption violations that were of concern. Can you do a better job of checking them, and if necessary, address them?

Examples:

- If batting averages are indeed heavy tailed or skewed, move from a normal distribution to something else. For example, could try a t-distribution with Laplace inference to handle the non-conjugacy.
- If the variance is indeed too high for a binomial model, try something that can handle the overdispersion.

## C Notes to self: some details for when slides are constructed

### C.1 Why Bayes

**Bayes law in real life** The primary motivation here is the importance of not “flipping the conditional” during interpretation. Sally Clark was convicted of murdering her two children, because the chance of two babies dying of SIDS in one family was one in 73m. But the expert witness ignored the prior probability that someone was a double murderer [7]. (See [14] for discussion.) Could also refer to the classic example with positive tests - maybe in reference to COVID [7]. *Why most published research findings are false* [8] provides a third example.

**Modeling Application: Estimating the probability of a rare event** See Section 1.2.1 of [2]. The problem is to estimate the proportion of people in a small town that have a disease given a small sample of 20 individuals. This is a nice believable example in which prior information is natural: we use prevalence in similar towns. (Make a note that really this is foreshadowing hierarchical models and hierarchical regressions!) This also illustrates another nice property of Bayes – we can get lots of functionals from the posterior. The plot on the right of Figure 1.1 – showing the shift between the prior and posterior distribution – is a classic. The sensitivity analysis is pretty cool – see the right hand plot of Figure 1.2, as well as the last couple sentences on pp.7.

The REAL kicker, I think, is the comparison to non-Bayesian methods. The frequentist confidence interval is complete garbage here, and the “adjusted” Wald interval is clearly just a (very specific) choice of prior. Nice opportunity for discussion here. Ask: what’s the advantage of Bayes over that? Possible answers: That is seemingly ad hoc; it’s not flexible to other choice of priors; Bayes makes it clear how it relates to priors; and it doesn’t allow for sensitivity analysis (or investigation of various functionals).

### C.2 Introduction to inference

**Bayes factors** Very nice brief discussion in Section 11.2.2 of [9].

**Multivariate normal - missing data and imputation** See Section 7.5 of [2].

### C.3 Hierarchical models

Some motivations:

- A way to use “surrounding data” as a prior in a more formal way. Think back to Hoff’s disease prevalence example. In that case, we constructed our beta prior manually, by taking a couple of basic facts about similar towns and then converting that into beta parameters. A hierarchical model could let the prior expectation be tied more exactly to those surrounding towns (including, if a regression is involved, similarity w.r.t. relevant characteristics, such as size, SES, etc.), to automatically set the strength of the prior expectation according to the relative uncertainty within and between towns, and to automatically adapt as data rolls in.
- Can be convenient. Consider Wand’s construction of the half-t distribution as a hierarchical model of inverse gammas, which makes for a conditionally conjugate scheme. Basically any auxiliary variable trick (think Polya Gamma augmentation) can be considered as the construction of a hierarchical model for computational convenience.
- By adding layers to the hierarchy, we can escape bad assumptions – e.g., the beta-binomial model handles over-dispersion. (This might be a better example for “why Bayes”.)