

Bayesian Statistics

Dustin Leatherman

October 16, 2020

Contents

1	Intro (2020/09/10)	3
1.1	Motivating Example	3
1.2	Frequentist Approach	4
1.2.1	Properties	4
1.2.2	Things a Frequentist would never say	4
1.3	Bayesian Approach	5
1.3.1	Likelihood Function	5
1.3.2	Priors	5
1.3.3	Posteriors	6
1.3.4	Advantages	6
1.3.5	Disadvantages	6
1.4	Review	7
1.4.1	Probability	7
1.4.2	Uncertainty	7
1.4.3	Probability vs Statistics	7
2	Probability & Introduction to Bayes (2020/09/17)	8
2.1	Calculating the Posterior Analytically	8
2.1.1	Using an Arbitrary PDF	8
2.1.2	Using Normal Distribution	10
2.2	Bayes Theorem	11
2.3	Bayesian Learning	12
2.4	Subjectivity	12
3	Summarizing a Posterior Distribution (2020/09/24)	12
3.1	SIR Model	12
3.2	Summarize a univariate Posterior with Beta-Binomial	13
3.3	MAP Estimator	13

3.4	Uncertainty Measures	13
3.4.1	Credible Sets	14
3.5	Hypothesis Tests	14
3.6	Monte Carlo Sampling	14
3.6.1	Transformations	14
3.7	Summarizing Multivariate Posteriors	15
3.8	Bayesian One Sample t-test	15
3.9	Frequentist vs Bayesian Analysis of a Normal Mean	16
3.10	Multiple Parameters in Multivariate Posteriors	16
3.11	Types of Uncertainty	17
3.11.1	Resolving Uncertainty	17
4	Conjugate and Objective Priors (2020/10/01)	18
4.1	Conjugate	18
4.1.1	Beta-Binomial	18
4.1.2	Related Problem using NegBin	21
4.1.3	Poisson-Gamma: One observation	21
4.1.4	Poisson-Gamma: Two Observations	22
4.1.5	Poisson-Gamma: m Observations	22
4.1.6	Gaussian-Gaussian	23
4.1.7	Gaussian-Gaussian: Known μ	24
4.2	Informative vs Uninformative	25
4.2.1	Mixture of Experts	26
4.3	Improper Priors	26
4.4	Subjective vs Objective Bayes	26
4.4.1	Objective Bayes	26
5	Deterministic Methods & MCMC Sampling (2020/10/08)	28
5.1	MAP Estimation (Maximum a Posteriori)	29
5.1.1	Example	29
5.2	Bayesian Central Limit Theorem	29
5.2.1	Example	30
5.3	Numerical Integration	31
5.4	Monte Carlo Sampling	31
5.4.1	Gibbs Sampling	32
6	MCMC Sampling & Convergence (2020/10/15)	34
6.1	Metropolis-Hastings Sampling	34
6.2	Gibbs Sampling	35
6.3	Metropolis Sampling	35

6.3.1	Tuning s_j	36
6.3.2	Logistic Regression Example	36
6.4	Convergence Diagnostics	36
6.4.1	Geweke	36
6.4.2	Gelman-Rubin	36
6.4.3	Effective Sample Size (ESS)	37
6.5	Handling Massive Datasets	37

1 Intro (2020/09/10)

1.1 Motivating Example

There are two students: Student A and Student B, along with an instructor. A secretly writes down a number (1,...,10) then mentally calls heads or tails.

1. The instructor flips a coin
2. If heads, A honestly tells B if the number is even or odd.
3. If A guesses H/T correctly, A tells B if their number is even or odd. Otherwise, they lie.
4. B will guess if the number is odd or even

Let θ be the probability that B correctly guesses even or odd.

The class (and myself) initially agreed without much discussion that 0.5 is the obvious answer. Upon further thinking on this, the probabilities breakdown in such way:

(2): 0.5 (3): 0.5 (4): ?

The initial logic is that its a 50/50 chance since there are two choices but there is an X-factor here with number 4. A few questions worth asking:

1. Does B know the rules upfront? As in, are they aware that A may or may not lie?
1. Does B see the result of the coin flip?
2. Is this done virtually or in person?

If the answer is no for 1 and 2, then 0.5 is a logical choice because they'd be guessing without much foreknowledge.

If the answer is yes for 1 and 2, then B is in on the "game" and can make a more educated guess. If A or the professor has a "tell", then that could provide information. Reading body language may also provide some information to B on the veracity of A's claim.

I would argue that θ would be > 0.5 if A and B know each other well enough. Which is really a great example of Bayesian vs Frequentist view points.

1.2 Frequentist Approach

Quantifies uncertainty in terms of repeating the process that generated the data many times.

1.2.1 Properties

- The parameters θ are fixed, unknown, and a constant.
- The sample (data) Y are random
- All prob. statements would be made about the randomness in the data.
-

A statistic $\hat{\theta} = Y/n$ is a statistic and is an estimator of the population proportion θ

The distribution of $\hat{\theta}$ from repeated sampling is the *sample distribution*.

1.2.2 Things a Frequentist would never say

- $P(\theta > 0) = 0.6$ because θ is not a random variable
- The distribution of θ is Normal(4.2,1.2)
- The probability that the true proportion is in the interval (0.4, 0.5) is 0.95.
- The probability that the null hypothesis is true is 0.03.

1.3 Bayesian Approach

Expresses uncertainty about θ using probability distributions. θ is still fixed and unknown.

Distribution *before* observing the data is the **prior distribution**. e.g. $P(\theta > 0.5) = 0.6$. This is subjective since people may have different priors.

Hopefully, Uncertainty about θ is reduced after observing the data.

Bayesian Interpretations differ from *Frequentist* Interpretations.

Uncertainty distribution of θ after observing the data is the **posterior distribution**.

Bayes Theorem for updating the prior

$$f(\theta|Y) = \frac{f(Y|\theta)f(\theta)}{f(Y)} \quad (1)$$

Described in words: Posterior \propto Likelihood \times Prior

$f(\theta|Y)$ is the posterior distribution.

Given that I have seen some data, what am I seeing now?

A key difference between Bayesian and frequentist statistics is that all inference is conditional on the single data set we observed (Y).

1.3.1 Likelihood Function

Distribution of the observed data given the parameters. This is the Same function used in a maximum likelihood analysis.

When prior information is weak, Bayesian and Maximum Likelihood Estimates are similar.

1.3.2 Priors

Say we observed $Y = 60$ successes in $n = 100$ trials and $\theta \in [0, 1]$ is the true probability of success.

Want to select a prior that has a domain of $[0, 1]$

If there is no relevant prior information, we might use $\theta \sim Uni(0, 1)$. This is called an *uninformative prior*. aka a “best guess”.

1. Beta

Beta distributions are a common prior for parameters between 0 and 1.

If $\theta \sim \text{Beta}(a, b)$, then the posterior is

$$\theta|Y \sim \text{Beta}(Y + a, n - Y + b)$$

$$\text{Beta}(1, 1) == \text{Uni}(0, 1)$$

2. Gamma Popular distribution for σ (population standard deviation)

1.3.3 Posteriors

The likelihood function $Y|\theta \sim \text{Bin}(n, \theta)$

The Uniform prior is $\theta \sim \text{Uni}(0, 1)$

The posterior is then $\theta|Y \sim \text{Beta}(Y + 1, n - Y + 1)$

1.3.4 Advantages

- Bayesian concepts (posterior probability of the null hypothesis) are arguably easier to interpret than the frequentist ideas (p-value.)
- Can incorporate scientific knowledge via the prior.
 - Even a Small amount of prior information can add stability.
- Excellent at quantifying uncertainty in complex problems.
- Provides a framework to incorporate data/information from multiple sources.

1.3.5 Disadvantages

- Less common/familiar
- Picking a prior is subjective (though there are objective priors)
- Procedures with frequentist properties are desirable.
- Computing can be slow for hard problems
- Non parametric methods are challenging

1.4 Review

Only the interesting parts are placed here. See the rest of this repo for deeper dives on other concepts.

1.4.1 Probability

Objective (associated with Frequentist)

- $P(X = x)$ is a mathematical statement
- If we repeatedly sampled X , the value that the proportion of draws equal to x converges is defined as $P(X = x)$

Subjective (associated with Bayesian)

- $P(X = x)$ represents an individual's degree of belief
- Often quantified as the amount an individual would be willing to wager that X will be x .

A Bayesian Analysis uses both of these concepts.

1.4.2 Uncertainty

Aleatoric (def: indeterminate) uncertainty (likelihood)

- Uncontrollable randomness in the experiment

Epistemic (def: involving knowledge) uncertainty (prior/posterior)

- Uncertainty about a quantity that could be theoretically

A Bayesian Analysis uses both of these concepts

1.4.3 Probability vs Statistics

The common sense, I like the way this is phrased.

Probability is the forward problems

- We assume we know how the data are being generated and compute the probability of events.

For example, what is the probability of flipping 5 straight heads if the coins are fair?

Statistics is the inverse problem

- We use data to learn about the data-generating mechanism

For example, if we flipped five straight heads, can we conclude the coin is biased?

2 Probability & Introduction to Bayes (2020/09/17)

if x and y are independent, then the following is true

$$f(x|y) = \frac{f(x,y)}{f_Y(y)} = \frac{f_X(x)f_Y(y)}{f_Y(y)} = f_X(x)$$

Cannot use $f(x,y)$ as PMF because $\sum_1^Y f(x,y) = f(x) \neq 1$. Need to scale by marginal probability in order to sum to 1 and thus be a proper PMF/PDF.

	1	2	3	4	5	Total [p(y)]
US	.0972	.0903	.0694	.0069	.0069	.2708
Not US	.3194	.1319	.1389	.1181	.0208	.7292
Total [p(x)]	.4167	.2222	.2083	.1250	.0278	1

show that x and y are dependent

$$P(x = 1) = 0.4167$$

$$P(y = 1) = 0.2708$$

$$P(x = 1) \times P(y = 1) = 0.4167(0.2708) = 0.1128$$

$$P(x = 1, y = 1) = 0.0972 \neq 0.1128 \text{ so dependent!}$$

2.1 Calculating the Posterior Analytically

2.1.1 Using an Arbitrary PDF

1. Find Joint Probability ($f(x,y)$)

$$\begin{aligned}
P(x > 7, y > 40) &= \int_7^{10} \int_{40}^{50} 0.26 \exp(-|x-7| - |y-40|) \, dx \, dy \\
&= 0.26 \int_7^{10} \int_{40}^{50} \exp(-x+7-y+40) \, dx \, dy \quad (\text{Since only interested in positive values}) \\
&= 0.26 \int_7^{10} \int_{40}^{50} \exp(-(x-7)) \exp(-(y-40)) \, dx \, dy \\
&= 0.26 \int_7^{10} \int_0^{10} \exp(-(x-7)) \exp(-u) \, dx \, du \\
&= 0.26 \int_7^{10} \int_0^{10} \exp(-(x-7)) [-\exp(-u)]_0^{10} \, dx \, du \\
&= 0.26(1 - e^{-10}) \int_7^{10} \exp(-(x-7)) \, dx \\
&= 0.26(1 - e^{-10})(1 - e^{-3}) \approx 0.247
\end{aligned} \tag{2}$$

1. Find Marginal Probability over the Data $f_X(x)$

$$\begin{aligned}
f_X(x) &= \int_{20}^{50} 0.26 + e^{-|x-7|-|y-40|} \, dy \\
&= 0.26 e^{-|x-7|} \int_{20}^{50} e^{-|y-40|} \, dy \\
&= 0.26 e^{-|x-7|} \left[\int_{20}^{40} e^{-(40-y)} \, dy + \int_{40}^{50} e^{-(y-40)} \, dy \right] \\
&= 0.26 e^{-|x-7|} \left[\int_{20}^0 -e^{-u} \, du + \int_0^{10} e^{-u} \, du \right] \\
&= 0.26 e^{-|x-7|} [1 - e^{-20} + 1 - e^{-10} \approx 2] \\
&= 0.52 e^{-|x-7|} \quad \forall x \leq x \leq 10
\end{aligned} \tag{3}$$

1. Calculate Conditional Probability

$$f(y|x) = \frac{f(x, y)}{f_X(x)} = \frac{1}{2} e^{-|y-40|}$$

If integrating over an absolute value, break up the integral into two integrals: the first over the negative domain of the integration, the second over the positive domain.

2.1.2 Using Normal Distribution

1. Find Marginal Probability

$$\begin{aligned}
 f(x) &= \int_{-\infty}^{\infty} \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{x^2 + y^2 - 2\rho xy}{2(1-\rho^2)}\right) dy \\
 &= \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-x^2/2(1-\rho^2)} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2 - 2\rho xy}{2(1-\rho^2)}\right) dy \quad (\text{Move x's out of integral. Arrange term}) \\
 &= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} e^{-x^2/2(1-\rho^2)} \int_{-\infty}^{\infty} \frac{\sqrt{1-\rho^2}}{\sqrt{2\pi}(1-\rho^2)} \exp\left(-\frac{y^2 - 2\rho xy + \rho^2 x^2 - (\rho x)^2}{2(1-\rho^2)}\right) dy \\
 &= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} e^{-x^2/2(1-\rho^2)} e^{\frac{\rho x^2}{2(1-\rho^2)}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}(1-\rho^2)} \exp\left(-\frac{(y - \rho x)^2}{2(1-\rho^2)}\right) dy, \quad N(\rho x, 1 - \rho^2) \\
 &= \frac{1}{\sqrt{2\pi}} e^{-0.5 \frac{x^2 - \rho^2 x^2}{1-\rho^2}} \\
 &= \frac{1}{\sqrt{2\pi}} e^{-0.5x^2}, X \sim N(0, 1)
 \end{aligned} \tag{4}$$

1. Assume Joint Normal PDF
2. Find Conditional probability

$$\begin{aligned}
f(y|x) &= \frac{f(x, y)}{f_X(x)} \\
&= \frac{\frac{1}{2\pi\sqrt{1-\rho^2}} \exp(-\frac{x^2+y^2-2\rho xy}{2(1-\rho^2)})}{\frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})} \\
&= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp(-\frac{x^2+y^2-2\rho xy}{2(1-\rho^2)} + \frac{x^2}{2}) \\
&= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp(-\frac{1}{2}[\frac{x^2+y^2-2\rho xy}{1-\rho^2} - x^2]) \tag{5} \\
&= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp(-\frac{1}{2}[\frac{x^2+y^2-2\rho xy - (1-\rho^2)x^2}{1-\rho^2}]) \\
&= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp(-\frac{1}{2}[\frac{y^2-2\rho xy - \rho^2 x^2}{1-\rho^2}]) \\
&= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp(-\frac{1}{2}[\frac{(y-\rho x)^2}{1-\rho^2}]), \quad y|x \sim N(\rho x, 1-\rho^2)
\end{aligned}$$

2.2 Bayes Theorem

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}$$

How do you know you are using Bayes Rule?

Given $P(y|\theta)$, want to find $P(\theta|y)$

- Bayesians quantify uncertainty about fixed but unknown parameters by treating

them as random variables.

- This requires that we set a prior distribution $\pi(\theta)$ to summarize uncertainty before observing the data.
- The distribution of the observed data given the model parameters is the *likelihood function*, $f(Y|\theta)$
 - The likelihood function is the most important piece of a Bayesian Analysis because it links the data and the parameters.

2.3 Bayesian Learning

The posterior distribution $P(\theta|Y)$ summarizes uncertainty about the parameters given the prior and data.

Reduction in uncertainty from prior to posterior represents **Bayesian Learning**

Bayes Theorem (again):

$$P(\theta|Y) = \frac{f(Y|\theta)\pi(\theta)}{m(Y)}$$

$m(Y) = \int F(Y|\theta)\pi(\theta)d\theta$: marginal distribution of the data and can usually be ignored.

2.4 Subjectivity

Choosing Likelihood function and a prior distribution are subjective.

If readers disagree with assumptions, findings will be rejected so assumptions must be justified theoretically and empirically.

3 Summarizing a Posterior Distribution (2020/09/24)

3.1 SIR Model

Susceptible-Infected-Recovered

At time t , $S_t + I_t + R_t = N$ where N is the population.

States evolved according to the following differential equations

$$\begin{aligned}\frac{dS_t}{dt} &= -\beta \frac{S_t I_t}{N} \\ \frac{dI_t}{dt} &= \beta \frac{S_t I_t}{N} - \Gamma I_t\end{aligned}\tag{6}$$

β : Controls rate of new infections

Γ : Controls recovery rate

We will use a discrete approx to these curves with hourly time steps.

So? $dt = \frac{1}{24}$

Goal: Fit SIR Model for given values of β and Γ

3.2 Summarize a univariate Posterior with Beta-Binomial

Posterior = Likelihood \times Prior

Say there is a parameter θ

Likelihood: $Y|\theta \sim \text{Bin}(N, \theta)$

Prior: $\theta \sim \text{Uni}(0, 1) \equiv \text{Beta}(1, 1)$

Posterior: $\theta|Y \sim \text{Beta}(Y + a, N - y + b)$

Peak of the Posterior is the MLE of the Likelihood function when using an uninformative prior.

3.3 MAP Estimator

Posterior Mode is call the max a posteriori (MAP) estimator

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(\theta|y) = \underset{\theta}{\operatorname{argmax}} \log[f(Y|\theta)] + \log[\pi(\theta)] \quad (7)$$

if prior is uniform, MAP is MLE assuming $Y|\theta \sim \text{Bin}(\theta, n)$.

3.4 Uncertainty Measures

Posterior Std. Dev. is one measure of uncertainty

- If approx Gaussian, can use empirical rule
- Analogous but fundamentally different than std error.
 - Std err is the standard deviation of $\hat{\theta}$'s sampling distribution

Do not call them call them **confidence** intervals. Called **Credible** Intervals in Bayesian Statistics.

Interval (l, u) is $100(1 - \alpha)\%$ posterior credible interval if $P(l < \theta < u|Y) = 1 - \alpha$

Interpretation: "Given the data and the prior, the probability that θ is between l and u is 0.95."

Confidence interval interpretation:

With 95% Confidence, θ is between l and u.

A Bayesian Posterior is a distribution for $\theta|Y$ whereas the sampling distribution is for $\hat{\theta}$. While their expected values both represent the true mean, the sampling distribution is not a distribution of θ hence why “Confidence” is used when in the interpretation. The Bayesian Posterior is a distribution of θ so the posterior can be used for the interpretation.

3.4.1 Credible Sets

Not unique.

Let q_τ be the τ quantile of the posterior of the posterior such that $P(\theta < q_\tau|Y) = \tau$. Then $(q_{0.0}, q_{0.95})$, $(q_{0.01}, q_{0.96})$, etc. are all valid 95% credible sets.

Equal-Tailed intervals: $(q_{\alpha/2}, q_{1-\frac{\alpha}{2}})$

Highest posterior density interval searches for the smallest interval that contains the proper probability

3.5 Hypothesis Tests

Conducted by computing posterior prob of each hypothesis.

$$P(\theta < 0.5|Y) = \int_0^{0.5} P(\theta|Y)d\theta$$

analogous but different than a p-value.

p-value: Assuming the null hypothesis is true, the probability we got X or a value more extreme is Y.

Bayesian Hypothesis Test: Given the prior and the data, the probability the null hypothesis is true is Y.

3.6 Monte Carlo Sampling

A useful tool for summarizing a posterior.

In MC sampling, we draw S samples from the posterior;

$$\theta', \dots, \theta^{(s)} \sim P(\theta|Y)$$

and use these samples to approx the posterior.

3.6.1 Transformations

MC sampling facilitates studying the **transformations** of parameters.

For example, the odds corresponding to θ are $\gamma = \frac{\theta}{1-\theta}$

$$\gamma^{(1)} = \frac{\theta^{(1)}}{1 - \theta^{(1)}}, \dots, \gamma^{(S)} = \frac{\theta^{(S)}}{1 - \theta^{(S)}} \quad (8)$$

How to approximate the posterior mean and variance of γ ?
 Transform the odds and use the draws to approximate θ 's posterior!

3.7 Summarizing Multivariate Posteriors

Univariate posteriors captured by a simple plot. Not easy or impossible to do with multivariate posteriors.

Let $\theta = (\theta_1, \dots, \theta_p)$.

Ideally, we reduced to the univariate marginal posteriors. Then the same ideas for univariate models apply

$$P(\theta_1|Y) = \int \dots \int P(\theta_1, \dots, \theta_p|Y) d\theta_2, \dots, d\theta_p$$

Can use Monte Carlo sampling to estimate these integrals.

Need to confirm the above statement

3.8 Bayesian One Sample t-test

Likelihood: $Y_i|\mu, \sigma \sim N(\mu, \sigma^2)$ indep over $i = 1, \dots, n$ Priors: $\mu \sim N(\mu_0, \sigma_0^2)$
 independent of $\sigma^2 \sim InvGamma(a, b)$

Typically we are interested in marginal posterior because it accounts for uncertainty about σ^2

Marginal Posterior: $f(\mu|Y) = \int_0^\infty P(\mu, \sigma^2|Y) d\sigma^2$, $Y = (Y_1, \dots, Y_n)$

if σ is known, the posterior of $\mu|Y$ is Gaussian and 95% Credible Interval is $E(\mu|Y) \pm Z_{0.975} SD(\mu|Y)$

if σ is unknown, the marginal (over σ^2) posterior of μ is t with $\nu = n + 2a$ degrees of freedom.

$$E(\mu|Y) \pm t_{0.975} SD(\mu|Y)$$

$SD(\mu|Y)$: Standard Deviation

Can summarize results best in a table with Posterior Mean, Posterior SD, and 95% Credible Set.

3.9 Frequentist vs Bayesian Analysis of a Normal Mean

Frequentist

Estimate of the μ is \bar{Y} If σ is known, the 95% C.I. is: $\bar{Y} \pm z_{0.975} \frac{\sigma}{\sqrt{n}}$

If σ is unknown, the 95% C.I. is: $\bar{Y} \pm t_{0.975, n-1} \frac{s}{\sqrt{n}}$

where t is the quantile of a t-distribution.

Bayesian

Estimate of μ is its marginal posterior mean.

Interval estimate is 95% Credible Interval.

If σ is known, Posterior of $\mu|Y$ is Gaussian

$E(\mu|Y) \pm Z_{0.975} SD(\mu|Y)$

If σ is unknown, the marginal (over σ^2) posterior of μ is t with $\nu = n + 2a$ degrees of freedom.

$E(\mu|Y) \pm t_{0.975, \nu} SD(\mu|Y)$

3.10 Multiple Parameters in Multivariate Posteriors

Want to compute $P(\theta_2 > \theta_1 | Y_1, Y_2)$.

Monte Carlo sampling of the posteriors a key tool!

Model is:

$$\begin{aligned} Y_1 | \theta_1 &\sim \text{Bin}(N, \theta_1) \\ Y_2 | \theta_2 &\sim \text{Bin}(N, \theta_2) \\ \theta_1, \theta_2 &\sim \text{Beta}(1, 1) \end{aligned} \tag{9}$$

Marginal Posteriors both independent of each other.

$$\bullet \theta_1 | Y_1, Y_2 \sim \text{Beta}(Y_1 + 1, N - Y_1 + 1)$$

$$\bullet \theta_2 | Y_1, Y_2 \sim \text{Beta}(Y_2 + 1, N - Y_2 + 1)$$

N <- 10; Y1 <- 5; Y2 <- 8;

S <- 10000

theta1 <- rbeta(S, Y1 + 1, N - Y1 + 1) theta2 <- rbeta(S, Y2 + 1, N - Y2 + 1)

(Y1 + 1) / (N + 2)

mean(theta1)

mean(theta2 > theta1)

3.11 Types of Uncertainty

Sampling

Parametric: Uncertainty about my guesses of the distribution of the parameter

3.11.1 Resolving Uncertainty

1. Plugin approach

If $\hat{\theta}$ is an estimate, thus $Y^* \sim f(Y|\hat{\theta})$

For example, Let $\hat{\theta} = \frac{2}{10}$. Predict $P(Y > 0) = 1 - (1 - 0.2)^{10}$.

If $\hat{\theta}$ has small uncertainty, this is fine. Otherwise, this underestimated uncertainty in Y^*

2. Posterior Predictive Distribution (PPD)

For the sake of prediction, the parameters aren't of interest as the parameters are vehicles by which the data inform about the predictive model.

PPD averages over their posterior uncertainty which *accounts* for parametric uncertainty.

$$f(Y^*|Y) = \int f(Y^*|\theta)p(\theta|Y) d\theta$$

Input = data Output = prediction distribution

Given I've observed a certain amount of data Y, what is the distribution of the predictor values?

Monte Carlo sampling approximates the PPD.

- (a) Example

Let $\theta^{(1)}, \dots, \theta^{(S)}$ be samples from the posterior.

Let $Y^{*(s)} \sim f(Y|\theta^{(s)})$ where $Y^{*(s)}$ are samples from the PPD for each $\theta^{(s)}$

Posterior Predictive Mean \approx sample mean of the $Y^{*(s)}$

$P(Y^* > 0) \approx$ sample proportion of non-zero $Y^{*(s)}$

$Y < -2$; $n < 10$;

$A <- Y + 1$; $B <- N - Y + 1$

```

1-dbinom(0,10,0.2)
theta <- rbeta(100000,A,B) Ystar <- rbinom(100000,10,theta)
mean(Ystar>0)

```

4 Conjugate and Objective Priors (2020/10/01)

How do we choose priors? This is the most important step of a Bayesian Analysis.

Key Terms

- Conjugate vs Non-conjugate
- Informative vs Uninformative
- Proper vs Improper
- Subjective vs Objective

4.1 Conjugate

Def: Prior and Posterior Distribution are from the same parametric family. This is done through a pairing of the Likelihood Distribution and the Prior Distribution.

4.1.1 Beta-Binomial

Use Case: Estimating a Proportion!

- What is the probability of success for a new cancer treatment?
- What proportion of voters support a candidate?

Let $\theta \in [0, 1]$ be a proportion we are trying to estimate.

Likelihood: $Y|\theta \sim \text{Bin}(n, \theta)$

Prior: $\theta \sim \text{Beta}(a, b)$

a: Prior number of successes b: Prior number of failures

Posterior: $\theta|Y \sim \text{Beta}(Y + a, n - Y + b)$

1. Frequentist Approach

MLE: $\hat{\theta} = \frac{Y}{n}$

$\hat{\theta} \sim N(\theta, \frac{\theta(1-\theta)}{n})$ for large Y and $n - Y$

Rule of Thumb for large enough n for proportions: At least 10-15 failures and 10-15 successes depending on which text book you read.

This is slightly different than the magic number 30 which is considered large enough for the mean.

$$SE(\hat{\theta}) = \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

2. Proof

(a) Short way

The short proof uses “proportional to” (\propto) and hand waves the constants.

Posterior:

$$\begin{aligned} f(\theta|Y) &\propto f(Y|\theta)f(\theta) = \binom{n}{Y} \theta^Y (1-\theta)^{n-Y} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1} \\ &\propto \theta^Y (1-\theta)^{n-Y} \theta^{a-1} (1-\theta)^{b-1} \\ &= \theta^{Y+a-1} (1-\theta)^{n-Y+b-1} \text{ (Looks like a Beta PDF)} \\ &\therefore \theta|Y \sim \text{Beta}(Y+a, n-Y+b) \end{aligned} \tag{10}$$

(b) Long way

$$f(Y|\theta) = \frac{f(Y|\theta) \cdot f(\theta)}{f(Y)} = \frac{f(Y|\theta) \cdot f(\theta)}{\int_0^1 f(Y, \theta) d\theta} \tag{11}$$

Numerator:

$$\begin{aligned} f(Y|\theta)f(\theta) &= \binom{n}{Y} \theta^Y (1-\theta)^{n-Y} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1} \\ &= \binom{n}{Y} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{Y+a-1} (1-\theta)^{n-Y+b-1} \end{aligned} \tag{12}$$

Denominator:

$$\begin{aligned}
f(Y) &= \int_0^1 \binom{n}{Y} \theta^Y (1-\theta)^{n-Y} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1} d\theta \\
&= \binom{n}{Y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \theta^{Y+a-1} (1-\theta)^{n-Y+b-1} d\theta \\
&= \binom{n}{Y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(Y+a)\Gamma(n-Y+b)}{\Gamma(n+a+b)} \int_0^1 \frac{\Gamma(n+a+b)}{\Gamma(Y+a)\Gamma(n-Y+b)} \theta^{Y+a-1} (1-\theta)^{n-Y+b-1} d\theta \\
\text{So? } f(y) &= \binom{n}{Y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(Y+a)\Gamma(n-Y+b)}{\Gamma(n+a+b)}
\end{aligned} \tag{13}$$

Posterior:

$$\begin{aligned}
f(\theta|Y) &= \frac{\binom{n}{Y} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{Y+a-1} (1-\theta)^{n-Y+b-1}}{\binom{n}{Y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(Y+a)\Gamma(n-Y+b)}{\Gamma(n+a+b)}} \\
&= \frac{\Gamma(n+a+b)}{\Gamma(Y+a)\Gamma(n-Y+b)} \theta^{Y+a-1} (1-\theta)^{n-Y+b-1} \\
\theta &\sim \text{Beta}(Y+a, n-Y+b)
\end{aligned} \tag{14}$$

3. Shrinkage

Posterior mean: $\hat{\theta}_B = E(\theta|Y) = \frac{Y+a}{n+a+b}$

Posterior mean is between the sample proportion ($\frac{Y}{n}$) and the prior mean: $\frac{a}{a+b}$

$$\hat{\theta}_B = w \frac{Y}{n} + (1-w) \frac{a}{a+b}$$

where $w = \frac{n}{n+a+b}$

- When n is large, $\hat{\theta}_B$ is closer to $\frac{Y}{n}$.
- as a and b grow, posterior mean more dependent on the prior.

Definition: The gravitation between the Likelihood function and the prior data. If there is

What prior to select if research show θ is between 0.6 and 0.8? $a = 7$, $b = 3$ because $\frac{7}{7+3} = 0.7$

4.1.2 Related Problem using NegBin

Estimate the number of successes (Y) before n failures.

θ : probability of success

$\theta \sim \text{Beta}(a, b)$

$Y|\theta \sim \text{NegBin}(n, \theta)$

$$\begin{aligned} f(\theta|Y) &\propto \binom{Y+n-1}{Y} \theta^Y (1-\theta)^n \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1} \\ &\propto \theta^Y (1-\theta)^n \theta^{a-1} (1-\theta)^{b-1} \\ &= \theta^{Y+a-1} (1-\theta)^{n+b-1} \\ &\sim \text{Beta}(Y+a, n+b) \end{aligned} \tag{15}$$

4.1.3 Poisson-Gamma: One observation

Goal: Estimate a rate!

Let $\lambda > 0$ be the rate to be estimated.

- Observations made over a period of N and observe $Y \in \{0, 1, 2, \dots\}$ events

- expected number of events: $N\lambda$

$\hat{\lambda} = \frac{Y}{n} = MLE$

Likelihood: $Y|\lambda \sim \text{Poisson}(N\lambda)$

Prior: $\lambda \sim \text{Gamma}(a, b)$

λ is continuous and positive so Gamma is a natural distribution to use for estimating the rate.

Posterior: $\lambda|Y \sim \text{Gamma}(a+Y, b+N)$

Interpretation

a: Prior number of events b: Prior observation time

1. Proof (Short Way)

$$\begin{aligned} f(\lambda|Y) &\propto \frac{e^{-N\lambda} (N\lambda)^Y}{Y!} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \\ &\propto e^{-N\lambda} \lambda^Y \lambda^{a-1} e^{-b\lambda} \\ &\propto e^{-(N+b)\lambda} \lambda^{Y+a-1} \text{ (Looks like a Gamma)} \\ &\therefore \lambda|Y \sim \text{Gamma}(Y+a, N+b) \end{aligned} \tag{16}$$

2. Shrinkage

The posterior mean is between the sample rate ($\frac{Y}{N}$) and the prior mean ($\frac{a}{b}$)

$$\hat{\lambda}_b = E(\lambda|Y) = \frac{Y+a}{N+b}$$

$$\hat{\lambda}_B = w \frac{Y}{N} + (1-w) \frac{Y+a}{N+b}$$

where $w = \frac{N}{N+b}$

What if we have no information about λ ? In general, make PDF Wide

What if λ is likely between 0.6 and 0.8? $a = 7$, $b = 10$ because $E(Y) = \frac{a}{b} = \frac{7}{10} = 0.7$

4.1.4 Poisson-Gamma: Two Observations

Likelihood:

$$\begin{aligned} f(Y_1, Y_2 | \lambda) &= f(Y_1 | \lambda) \cdot f(Y_2 | \lambda) \text{ (if Y's are independent)} \\ &= \frac{(N\lambda)^{Y_1} e^{-N\lambda}}{Y_1!} \cdot \frac{(N\lambda)^{Y_2} e^{-N\lambda}}{Y_2!} \\ &= \frac{(N\lambda)^{Y_1+Y_2} e^{-2N\lambda}}{Y_1! \cdot Y_2!} \end{aligned} \quad (17)$$

Prior: $f(\lambda) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}$

Posterior:

$$\begin{aligned} f(\lambda | Y_1, Y_2) &\propto \frac{(N\lambda)^{Y_1+Y_2} e^{-2N\lambda}}{Y_1! Y_2!} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \\ &\propto \lambda^{Y_1+Y_2} e^{-2N\lambda} \lambda^{a-1} e^{-b\lambda} \\ &= \lambda^{Y_1+Y_2+a-1} e^{-2N\lambda-b\lambda} \text{ (Looks like a Gamma)} \\ &\therefore \lambda | Y_1, Y_2 \sim \text{Gamma}(Y_1 + Y_2 + a, 2N + b) \end{aligned} \quad (18)$$

4.1.5 Poisson-Gamma: m Observations

Likelihood:

$$\begin{aligned}
f(Y_1, \dots, Y_m | \lambda) &= f(Y_1 | \lambda) \cdot \dots \cdot f(Y_m | \lambda) \text{ (if Y's are independent)} \\
&= \prod_1^m \frac{(N\lambda)^{Y_i} e^{-N\lambda}}{Y_i} \\
&= \frac{(N\lambda)^{\sum Y_i} e^{-mN\lambda}}{\prod_1^m Y_i!}
\end{aligned} \tag{19}$$

Prior: $f(\lambda) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}$
Posterior:

$$\begin{aligned}
f(\lambda | Y_1, \dots, Y_m) &\propto \frac{(N\lambda)^{\sum Y_i} e^{-2mN\lambda}}{\prod_1^m Y_i! \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}} \\
&\propto (N\lambda)^{\sum Y_i} e^{-mN\lambda} \lambda^{a-1} e^{-b\lambda} \\
&\propto (N\lambda)^{a-1+\sum Y_i} e^{-(mN+b)\lambda} \text{ (Looks like a Gamma PDF)} \\
\therefore \lambda | Y_1, \dots, Y_m &\sim \text{Gamma}\left(\sum_1^m Y_i + a, mN + b\right)
\end{aligned} \tag{20}$$

4.1.6 Gaussian-Gaussian

Goal: Estimate a mean! (μ)

Likelihood: $f(Y_1, \dots, Y_n | \mu) = \frac{1}{\sigma^n (2\pi)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_1^n (Y_i - \mu)^2\right)$

Prior:

$$\begin{aligned}
f(\mu) &= \frac{1}{\sqrt{2\pi} \frac{\sigma}{\sqrt{m}}} \exp\left(-\frac{1}{2 \frac{\sigma^2}{m}} (\mu - \theta)^2\right) \\
&= \frac{\sqrt{m}}{\sqrt{2\pi}\sigma} \exp\left(-\frac{m}{2\sigma^2} (\mu - \theta)^2\right)
\end{aligned} \tag{21}$$

Posterior:

$$\begin{aligned}
f(\mu|Y_1, \dots, Y_n) &\propto \frac{1}{\sigma^n (2\pi)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum (Y_i - \mu)^2\right) \cdot \frac{\sqrt{m}}{\sqrt{2\pi}\sigma} \\
&\propto \exp\left(-\frac{1}{2\sigma^2} [\sum (Y_i - \mu)^2 + m(\mu - \theta)^2]\right) \\
&\propto \exp\left(-\frac{1}{2\sigma^2} [\sum (Y_i^2 - 2\mu Y_i + \mu^2) + m(\mu^2 - 2\mu\theta + \theta^2)]\right) \\
&\propto \exp\left(-\frac{1}{2\sigma^2} [2\mu \sum Y_i + n\mu^2 + m\mu^2 - 2m\mu\theta]\right) \text{ (where does the square } Y_i \text{ go?)} \\
&\propto \exp\left(-\frac{1}{2\sigma^2} [2\mu n\bar{Y} + n\mu^2 + m\mu^2 - 2m\mu\theta]\right) \\
&= \exp\left(-\frac{n+m}{2\sigma^2} \left[-2\frac{n\bar{Y} + m\theta}{n+m} + \mu^2\right]\right) \\
&\propto \exp\left(-\frac{n+m}{2\frac{\sigma^2}{n+m}} \left[\mu - \frac{n\bar{Y} + m\theta}{n+m}\right]^2\right) \text{ (Looks like a Normal PDF)} \\
&\therefore \mu|Y_1, \dots, Y_n \sim N\left(\frac{n\bar{Y} + m\theta}{n+m}, \frac{\sigma^2}{n+m}\right)
\end{aligned} \tag{22}$$

This can also be written as

Let $w = \frac{n}{n+m}$, then $\mu|Y_1, \dots, Y_m \sim N(w\bar{Y} + (1-w)\theta, \frac{\sigma^2}{n+m})$

m can *loosely* be interpreted as the prior number of observations

1. Shrinkage

$$\hat{\mu}_B = E(\mu|Y_1, \dots, Y_n) = w\bar{Y} + (1-w)\theta \text{ where } w = \frac{n}{n+m}$$

If no prior information available, make m small to make the prior uninformative. This is because a small m makes the variance large which makes the bell curve wide.

4.1.7 Gaussian-Gaussian: Known μ

If μ is known, then we should be estimating σ^2 .

$\sigma^2 \sim \text{Gamma}(a, b)$ since Gamma is continuous over $(0, \infty)$ which matches the domain of the variance.

The math is easier if using a gamma prior for the inverse variance (τ).

Inverse Variance is also known as *precision* $\frac{1}{\sigma^2}$

If $\frac{1}{\sigma^2} \sim \text{Gamma}(a, b)$, then $\sigma^2 \sim \text{InvGamma}(a, b)$

Likelihood: $f(Y_1, \dots, Y_n|\sigma^2) = \frac{1}{(\sqrt{2\pi})^n (\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_n (Y_i - \mu)^2\right)$

Prior: $f(\sigma^2) = \frac{b^a(\sigma^2)^{-a-1}e^{-b/\sigma^2}}{\Gamma(a)}$

Posterior:

$$\begin{aligned}
f(\sigma^2|Y_1, \dots, Y_n) &\propto \frac{1}{\sqrt{2\pi}^n (\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_1^n (Y_i - \mu)^2\right) \cdot \frac{b^a(\sigma^2)^{-b-1}}{\Gamma(a)} \exp\left(\frac{-b}{\sigma^2}\right) \\
&\propto (\sigma^2)^{-n/2} \exp\left(\frac{-1}{2\sigma^2} \sum_1^n (Y_i - \mu)^2\right) \\
&\propto (\sigma^2)^{-(\frac{n}{2}+a)-1} \exp\left(\frac{-1}{\sigma^2} \left[\frac{\sum_1^n (Y_i - \mu)^2}{2} + b\right]\right)
\end{aligned} \tag{23}$$

if μ is known, then $SSE = \sum_1^n (Y - \mu)^2$

$\sigma^2|Y_1, \dots, Y_n \sim InvGamma(\frac{n}{2} + a, \frac{SSE}{2} + b)$

Using τ

If Y_i has mean μ and precision τ , then likelihood is proportional to:

$$\Pi_n^1 f(y_i|\mu) \propto \tau^{n/2} \exp\left(-\frac{\tau}{2} \sum_1^n (y_i - \mu)^2\right)$$

If $\tau \sim Gamma(a, b)$, then $\tau|Y \sim Gamma(\frac{n}{2} + a, \frac{SSE}{2} + b)$

This matches the results when using an Envisages for Variance.

1. Shrinkage

Mean of InvGamma only exists for $a > 1$.

Prior mean: $\frac{b}{a-1}$

Posterior mean: $\frac{SSE+b}{n+2a-2}$

common to make a,b small to give an uninformative prior. Then posterior mean converges towards sample variance.

4.2 Informative vs Uninformative

Can use informative priors from literature reviews, pilot studies, expert opinions, etc.

Prior Elicitation: Process of converting expert information *into* a prior. Experts may not know what an InvGamma is but their information can be converted into one!

Weak/Uninformative Priors commonplace. Easier to defend.

Strong Priors typically used for nuisance parameters. i.e. parameters we don't care about. The idea being that we don't care about it and really only want to affect the analysis if its *really* strong.

Sensitivity Analysis used to compare the posterior against several priors. This lets readers know how the prior exactly affects the analysis.

4.2.1 Mixture of Experts

Combine several priors into a single prior. For example, there are three studies which promote three different priors. These can be combined into a single prior.

$$\pi(\theta) = \sum_{j=1}^J w_j \pi_j(\theta)$$

where w_j is a weight with the constraints $w_j > 0$ and $\sum w_j = 1$

4.3 Improper Priors

A prior that doesn't integrate to 1. e.g. $\pi(\mu) = 1 \forall \mu \in \mathbb{R}$

It is okay to use an improper prior as long as you verify the posterior integrates to 1.

4.4 Subjective vs Objective Bayes

An objective analysis requires no subjective decisions by the analyst such as picking a prior, picking a likelihood function, treatment of outliers or transforms, etc.

Objective analysis may be feasible in a tightly controlled study but generally impossible for most analysis.

4.4.1 Objective Bayes

Lets an algorithm choose prior.

Examples

- Jeffrey's Prior
- Probability matching
- Maximum Entropy
- Empirical Bayes

- Penalized Complexity

Jeffrey's priors are most common.

Most of these are *improper* so posterior needs to be checked that it integrates to 1.

1. Jeffrey's Prior

Jeffrey's Prior for θ : $p(\theta) = \sqrt{I(\theta)}$ where $I(\theta)$ is the Fisher Information Matrix.

$$I(\theta) = -E_{Y|\theta}[\frac{d^2}{d\theta^2} \log p(Y|\theta)]$$

Once the likelihood is specified, Jefferey's prior is determined with no additional input hence being objective about prior.

Examples

Likelihood: $Y \sim \text{Bin}(n, \theta)$

Jefferey's Prior: $\theta \sim \text{Beta}(0.5, 0.5)$

Likelihood: $Y \sim N(\mu, 1)$

Jefferey's Prior: $p(\mu) \propto 1$

Likelihood: $Y \sim N(0, \sigma^2)$

Jefferey's Prior: $p(\sigma) \propto 1/\sigma$

2. Reference Priors

Try to be uninformative. Univariate models give Jeffreys Priors. Multivariate models give different priors. Harder to compute than Jeffrey's.

3. Probability Matching Priors (PMP)

Designed so Posterior Credible Intervals have correct frequentist coverage.

For example, if $Y_i|\mu \sim N(\mu, 1)$, the PMP is $p(\mu) = 1$. Then, Posterior is $\mu|Y \sim N(\bar{Y}, 1/n)$

Only a few cases where this can be used.

4. Empirical Bayes

Pick priors based on data.

Ex: Maybe σ^2 has prior mean s^2

Criticized for using data twice: once for the prior, and once for the likelihood.

5. Penalized Complexity Priors (PCP)

A PCP prior begins with a simple base model. e.g linear regression with all slopes equal to 0.

Full model is shrunk towards base model. e.g regression with non-zero slopes.

Distance from full to base model has exponential prior to penalize the more complex model from deviating from the base.

Requires picking the parameter in the exponential prior and setting priors for the parameters in the base model **so not purely objective**

6. Maximum Entropy Priors

Entropy is a measure of uncertainty. The entropy of the PMF $f(x)$ is

$$-\sum_{x \in S} f(x) \log[f(x)]$$

- (a) Fix a few quantities of the prior distribution. e.g. $E(\theta) = 0.5$
- (b) Find the prior with maximum entropy that satisfies these constraints.

If θ has support \mathbb{R} and mean and variance are known, maximum entropy prior is Gaussian.

Not purely objective because you have to set the constraints

*

5 Deterministic Methods & MCMC Sampling (2020/10/08)

The big question is **How to summarize the Posterior?**

We need point estimates, credible sets, etc.

Algorithms to Estimate Complicated Joint Posteriors

- Use a point estimate (e.g. MAP), ignore uncertainty
- Approximate Posterior as Gaussian using Bayesian Central Limit Theorem
- Numerical Integration. Not touched on much here. Moreso in Numerical Analysis
- Markov-Chain Monte Carlo Sampling

5.1 MAP Estimation (Maximum a Posteriori)

Sometimes you don't need an entire posterior distribution. A single point estimate will do. For example, prediction in Machine Learning.

MAP Estimate is the posterior **mode**. AKA the peak of the posterior distribution.

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta|Y) = \underset{\theta}{\operatorname{argmax}} \log[f(Y|\theta)] + \log(\pi(\theta))$$

Similar to MLE but includes prior.

5.1.1 Example

Let $Y|\theta \sim \text{Gamma}(Y + a, N + b)$ and $\theta \sim \text{Gamma}(a, b)$. Find $\hat{\theta}_{MAP}$

$$p(\theta|Y) = \frac{(N + b)}{Y + a} \Gamma(Y + a) \theta^{Y+a-1} e^{-(N+b)\theta}$$

$$\begin{aligned} \hat{\theta}_{MAP} &= \underset{\theta}{\operatorname{argmax}} \log[f(\theta|Y)] \\ &= \underset{\theta}{\operatorname{argmax}} \{ (Y + a) \log(N + b) - \log(\Gamma(Y + a)) + (Y + a - 1) \log(\theta) - (N + b) \theta \} \\ &= \frac{d}{d\theta} \log(p(\theta|Y)) = \frac{y + a - 1}{\theta} - (N + b) = 0 \\ &= \frac{Y + a - 1}{N + b} \end{aligned} \tag{24}$$

5.2 Bayesian Central Limit Theorem

Berstein-Von Mises Theorem: As the sample size grows, the posterior doesn't depend on the prior. i.e. Shrinkage.

Def: For large N and some other conditions, $\theta|Y \approx \text{Normal}$

$$\theta|Y \sim N(\hat{\theta}_{MAP}, I(\hat{\theta}_{MAP})^{-1})$$

I is Fisher's Information Matrix, aka the Hessian Matrix.

$$I_{jk} = \frac{-d^2}{d\theta_j d\theta_k} \log[p(\theta|Y)]$$

evaluated at $\hat{\theta}_{MAP}$

5.2.1 Example

Let $\theta \sim \text{Beta}(0.5, 0.5)$ and $Y|\theta \sim \text{Bin}(n, \theta)$. Find Gaussian approximation for $p(\theta|Y)$.

In this case, $\text{Beta}(0.5, 0.5)$ is Jeffreys Prior.

Posterior: $\theta|Y \sim \text{Beta}(Y + 0.5, n - Y + 0.5)$

Need a MAP Estimator.

$$\frac{\Gamma(Y + 1 + n - Y)}{\Gamma(Y + 0.5)\Gamma(n - Y + 0.5)} \Theta^{Y-0.5}(1 - \theta)^{n-Y-0.5}$$

$$\begin{aligned} \log p(\theta|Y) &= \log \Gamma(n + 1) - \log \Gamma(Y + 0.5) - \log \Gamma(n - y + 0.5) + (y - 0.5)\log \theta + (n - y - 0.5)\log \theta \\ \frac{d}{d\theta} \log p(\theta|Y) &= \frac{Y - 0.5}{\theta} - \frac{n - y - 0.5}{1 - \theta} = 0 \\ \Rightarrow \frac{Y - 0.5}{\theta} &= \frac{n - Y - 0.5}{1 - \theta} \\ \Rightarrow \hat{\theta}(n - Y - 0.5) &= (1 - \hat{\theta})(Y - 0.5) \\ \Rightarrow \hat{\theta}(n - Y - 0.5) + \hat{\theta}(Y - 0.5) &= Y - 0.5 - \hat{\theta}(Y - 0.5) \\ \Rightarrow \hat{\theta}(n - Y - 0.5) &= Y - 0.5 \\ \Rightarrow \hat{\theta}(n - Y - 0.5 + Y - 0.5) &= Y - 0.5 \\ \Rightarrow \hat{\theta}(n - 1) &= Y - 0.5 \\ \therefore \hat{\theta}_{MAP} &= \frac{Y - 0.5}{n - 1} \end{aligned} \tag{25}$$

Finding the Variance via the Information Matrix

$$\begin{aligned} \frac{-d^2}{d\theta^2} \log p(\theta|Y) &= \frac{-d}{d\theta} \left[\frac{d}{d\theta} \log p(\theta|Y) \right] \\ &= \frac{-d}{d\theta} \left[\frac{Y - 0.5}{\theta} - \frac{n - Y - 0.5}{1 - \theta} \right] \\ &= - \left[-\frac{Y - 0.5}{\theta^2} - \frac{n - Y - 0.5}{(1 - \theta)^2} \right] \\ &= \frac{Y - 0.5}{\theta^2} + \frac{n - Y - 0.5}{(1 - \theta)^2} \\ &= \frac{Y - 0.5}{\theta^2} + \frac{n - Y - 0.5}{[1 - (\frac{Y - 0.5}{n - 1})]^2} \end{aligned} \tag{26}$$

$$\begin{aligned}
I(\hat{\theta}_{MAP}) &= \frac{Y - 0.5}{[\frac{Y-0.5}{n-1}]^2} + \frac{n - Y - 0.5}{[1 - \frac{Y-0.5}{n-1}]^2} \\
&= \frac{(n-1)^2}{Y - 0.5} + \frac{n - Y - 0.5}{[\frac{n-1}{n-1} - \frac{Y-0.5}{n-1}]^2} \\
&= \frac{(n-1)^2}{Y - 0.5} + \frac{n - Y - 0.5}{[\frac{n-Y-0.5}{n-1}]^2} \\
&= \frac{(n-1)^2}{Y - 0.5} + \frac{(n-1)^2}{n - Y - 0.5} \\
&= (n-1)^2 \frac{1}{Y - 0.5} + \frac{1}{n - Y - 0.5} \\
&= (n-1)^2 \frac{n-1}{(Y - 0.5)(n - Y - 0.5)} \\
&= \frac{(n-1)^3}{(Y - 0.5)(n - Y - 0.5)}
\end{aligned} \tag{27}$$

$$\therefore \theta|Y \approx N\left(\frac{Y-0.5}{n-1}, \frac{(Y-0.5)(n-Y-0.5)}{(n-1)^3}\right)$$

Note that $I(\hat{\theta}_{MAP})$ produces the Inverse Variance hence why the Variance of the approximation is flipped.

Normal Approximation for Beta-Binomial using Bayesian Central Limit Theorem.

For large Datasets with small number of params, the Normal approximation is good.

5.3 Numerical Integration

Only feasible for small p.

Iteratively Nested Laplace Approximation (INLA) combines Gaussian approximation with numerical integration. It works well if most parameters are approximately normal.

5.4 Monte Carlo Sampling

Collection of all params in model: $\theta = (\theta_1, \dots, \theta_p)$

Dataset: $Y = (Y_1, \dots, Y_n)$

Posterior Distribution: $f(\theta|Y)$

If $\theta^{(1)}, \dots, \theta^{(s)}$ are samples from $f(\theta|Y)$, then mean of the S samples approximate a posterior mean.

Most common MCMC Algorithms: Gibbs, Metropolis

5.4.1 Gibbs Sampling

- Sample from High dimension Posteriors
- Break problem of sampling from the high dimension joint distribution into a series of samples from low dimensional conditional distributions.

Samples are not independent. The dependencies form a Markov Chain.

Full Conditional (FC) Distribution: Distribution of one parameter taking all other parameters as *fixed and known*.

1. MCMC for Bayesian T-Test

$Y_i \sim N(\mu, \sigma^2)$ where $\mu \sim N(0, \sigma_0^2)$ and $\sigma^2 \sim InvGamma(a, b)$

FC Distribution 1

$$\mu|\sigma^2, Y \sim N\left(\frac{n\bar{Y}\sigma^{-2} + \mu_0\sigma_0^{-2}}{n\sigma^{-2} + \sigma_0^{-2}}, \frac{1}{n\sigma^{-2} + \sigma_0^{-2}}\right)$$

FC Distribution 2

$$\sigma^2|\mu, Y \sim InvGamma\left(\frac{n}{2} + a, \frac{1}{2} \sum_{i=1}^n (Y_i - \mu)^2 + b\right)$$

2. Algorithm

- (a) Set Initial values for all parameters: $\theta_1^{(0)}, \dots, \theta_p^{(0)}$
- (b) Sample variables one at a time from Full Conditional: $p(\theta_j|\theta_1, \dots, \theta_{j-1}, \theta_{j+1}, \dots, \theta_p, Y)$
 - Rather than 1 p-dimensional samples, we produce p 1-dimensional samples
 - Repeat until required number of samples are generated.

Example

- (a) Set initial values
- (b) For iteration t

- FC1: Draw $\theta_1^{(t)} | \theta_2^{(t-1)}, \dots, \theta_p^{(t-1)}, Y$
- FC2: Draw $\theta_2^{(t)} | \theta_1^{(t)}, \theta_3^{(t-1)}, \dots, \theta_p^{(t-1)}, Y$
- ...
- FCp: Draw $\theta_p^{(t)} | \theta_1^{(t)}, \dots, \theta_{p-1}^{(t)}, Y$

Repeat 2 \mathbb{S} times giving posterior draws $\theta^{(1)}, \dots, \theta^{(S)}$

Why does it work?

Theorem: For any initial values, the chain will eventually converge to the posterior.

Theorem: If $\theta^{(s)}$ is a sample from the posterior, then $\theta^{(s+1)}$ is too.

Once the chain has *converged*, then discard the first T samples as “burn in”. Use remaining S - T to approximate the Posterior.

I have also heard “burn in” described as annealing in my Monte Carlo class.

3. Practice Problem

Work out the full conditionals for λ and b for the following model:

$$\begin{aligned} Y | \lambda, b &\sim \text{Poisson}(\lambda) \\ \lambda | b &\sim \text{Gamma}(1, b) \\ b &\sim \text{Gamma}(1, 1) \sim \text{exp}(1) \end{aligned} \tag{28}$$

Posterior:

$$\begin{aligned} P(\lambda, b | Y) &\propto f(Y | \lambda, b) f(\lambda | b) = f(Y | \lambda, b) f(b) f(\lambda | b) \\ &= \frac{e^{-\lambda} \lambda^Y}{Y!} \cdot e^{-b} \cdot b e^{-b\lambda} \\ &\propto e^{-\lambda} \lambda^Y e^{-b} b e^{-b\lambda} \\ &= b e^{-\lambda - b - b\lambda} \lambda^Y \end{aligned} \tag{29}$$

Full Conditional Distributions

- (a) Write in terms of λ . Everything else is fixed.

$$\begin{aligned}
\lambda|b, Y &\propto \lambda^Y e^{-\lambda} e^{-b\lambda} \\
&= \lambda^{(y+1)-1} e^{-(b+1)\lambda} \\
\text{So? } \lambda|b, Y &\sim \text{Gamma}(y+1, b+1)
\end{aligned} \tag{30}$$

(a) Write in terms of b . Everything else is fixed.

$$\begin{aligned}
b|\lambda, Y &\propto b e^{-b} e^{-b\lambda} = b e^{-(\lambda+1)b} \\
\text{So? } b|\lambda, Y &\sim \text{Gamma}(2, \lambda+1)
\end{aligned} \tag{31}$$

6 MCMC Sampling & Convergence (2020/10/15)

Adaptive MCMC: Uses same candidate distribution throughout the chain but adjusts hyper-parameters to tune Acceptance Rate.

Hamiltonian MCMC: Uses gradient of the posterior to adjust hyper-parameters throughout the chain. Default MCMC method in STAN.

3 Main Decisions for MCMC Algo

1. Select Initial Values
 - Use Method of Moments or MLE
 - Pick purposefully bad values to demonstrate convergence.
2. Determine the chain convergence
3. Determine how many samples you need (Effective Sample Size)

The following MCMC Sampling methods can be mixed and matched as needed to use and sample from candidate distributions (barring Gibbs).

6.1 Metropolis-Hastings Sampling

Generic case which allows for Asymmetric Candidate Distributions.

Let θ_j^c be a Random Variable from the candidate distribution

$$\theta_j^c \sim q(\theta|\theta^*)$$

For example, if $\theta_j^c \sim N(\theta_j^*, s_j^2)$, then

$$q(\theta_j^c | \theta^*) = \frac{2}{s_j \sqrt{2\pi}} \exp\left[-\frac{(\theta_j^c - \theta_j^*)^2}{2s_j^2}\right] = q(\theta^* | \theta_j^c)$$

High correlation cases cause slow convergence. If high correlation is present, **Blocked Gibbs/Metropolis** can help.

Ex. Linear Regression iterates between sampling the block $(\beta_1, \dots, \beta_p)$ and σ^2

Blocked still not clear. Follow up with more.

6.2 Gibbs Sampling

Gibbs samples each parameter from its conditional distribution. Which makes use of conjugate priors. It is not obvious to use without conjugate priors.

$$P(\mu | Y) \propto \exp\left[-\frac{1}{2}(Y - \mu)^2\right] \cdot \mu^{a-1}(1 - \mu)^{b-1}$$

Where $Y \sim N(\mu, 1)$, $\mu \sim \text{Beta}(a, b)$

No known conjugate prior for some likelihoods. e.g. Logistic Regression. This is where Metropolis Sampling comes in!

Special Case of Metropolis sampling where acceptance rate is always 1 and the proposal distributions are the posterior conditionals.

6.3 Metropolis Sampling

A version of rejection Sampling

Let θ_j^* be the current value of the parameter being updated and θ_j be the current value for all parameters.

Let θ_j^c be the candidate value where

$$\theta_j^c \sim N(\theta_j^*, s_j^2)$$

Let R be the probability of accepting a move where

$$R = \min\left\{1, \frac{P(\theta_j^c | \theta_{(j)}, Y)}{P(\theta_j^* | \theta_{(j)}, Y)}\right\}$$

This is a special case of Metropolis-Hastings sampling where the candidate distribution is **symmetric**.

6.3.1 Tuning s_j

s_j is a hyper-parameter for the Metropolis Algorithm. Ideally, it is somewhere between 0.3 and 0.4. This is because we don't want to accept or reject *too* much.

If s_j is small

- proposed distribution is narrow
- Nearly all candidates are accepted

If s_j is large

- proposed distribution is wide
- nearly all candidates rejected. As in, there are a lot of straight lines on the trend graph.

6.3.2 Logistic Regression Example

This is how Bayesians see Logistic Regression.

$$Y_i|\theta \sim \text{Bern}(\text{logit}^{-1}(\theta)) = \text{Bern}\left(\frac{e^\theta}{1 + e^\theta}\right)$$

Where $\theta \sim N(\mu_0, \sigma_0^2)$

6.4 Convergence Diagnostics

6.4.1 Geweke

Compares mean at beginning of chain with mean at end of chain.

Test statistic: Z-score.

$|Z| > 2 \implies$ poor convergence.

6.4.2 Gelman-Rubin

By running multiple chains, hopefully to see the same result. The measurements between chains should agree. Essentially an ANOVA test of whether the chains have the same mean.

1 is perfect 1.1 is decent but not great convergence.

6.4.3 Effective Sample Size (ESS)

Idea

Highly correlated samples have less information than independent samples. ESS accounts for this autocorrelation. For example, you may think you have 10,000 samples but you actually have less than 10K *good* samples because you need to account for autocorrelation.

S : Number of MCMC Samples

$\rho(h)$: Autocorrelation with Lag h .

For example, $\rho(1)$ is autocorrelation coefficient with the first Lag.

$$ESS = \frac{S}{1 + 2 \sum_{h=1}^{\infty} \rho(h)}$$

ESS should be at least a few **thousand**.

Naive SE: $SE = \frac{S}{\sqrt{S}}$

Time Series SE is more realistic.

$$SE = \frac{S}{\sqrt{ESS}}$$

6.5 Handling Massive Datasets

- MAP Estimate (Prediction only. Good for ML)
- Bayesian CLT
- Variation Bayes: Approximate Posterior by assuming posterior independent across all parameters. (Fast but questionable statistical properties).
- Parallel Computing
- Batch Datasets to Divide and Conquer

Misc thing about JAGS, STAN.

Sometimes people use a cauchy distribution to do a standard deviation in JAGS, STAN.