

BAYESIAN INFERENCE AND REGULARIZATION

1. Learning objectives

- Understand the motivation for and basic definition of Bayesian inference.
- Familiarity with Bayesian inference for Bernoulli random variables using the β -distribution as priors.
- Understand how we select priors for our models.
- Understand how to compute (in python) the posterior statistics of regression parameters and make posterior predictions.

2. Introduction to Bayesian inference: some simple examples

- Models encode our assumptions about how data was generated. Sometimes we want to incorporate vague information, such as “the function describing my data is very smooth and changes roughly on a time-scale of 5 hours”. Or, we might want to include many predictors, but to avoid overfitting penalize high values of the predictors. For example, we might believe there is an interaction term, but suspect it is much smaller than the additive terms. This motivates a different approach to statistics, one which treats parameters as themselves random variables.

Example 1 (Bernoulli with priors). *Consider the model*

$$X \sim \text{Bernoulli}(q)$$

An implicit assumption we have made when we fit this model, is that before we see any data, each value of q is equally likely. To make this more explicit, we can write

$$X|q \sim \text{Bernoulli}(q)$$

$$q \sim \text{Uniform}(0, 1)$$

Another way to think about constructing an estimator of q , is that given data points X_1, \dots, X_N , we would like to know

$$(1) \quad \hat{q}_b = E[q|X_1, \dots, X_N]$$

We can in principle solve this problem using Bayesian Theorem. We would also take it a step further and ask: What is the distribution of

$$q|X_1, \dots, X_N$$

It turns out that this has a relatively simple form, called a β -distribution. We say

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

if q is a random variable on $[0, 1]$ which has density

$$f(q) = Bq^{a-1}(1-q)^{b-1}$$

The constant B ensures the area under this curve between $q = 0$ and $q = 1$ is indeed one, as it must be for a random density.

Question: Show that

$$q|X_1, \dots, X_N \sim \text{Beta}(Y + 1, N - Y + 1)$$

Solution: Using Bayes' theorem

$$\begin{aligned} f(q|Y) &= \frac{P(Y|q)f(q)}{P(Y)} \\ &= \frac{1}{P(Y)} \binom{N}{Y} q^Y (1-q)^{N-Y} \end{aligned}$$

Since we are thinking of this as a probability density in q , we don't actually need to compute $P(Y)$ explicitly (even though we could), instead, we can simply notice that

$$f(q|Y) = Cq^Y(1-q)^{N-Y}$$

for some constant of proportionality C which will depend on Y , but is uniquely determined by the fact that the area under this curve must be 1. In the colab notebook we plot this curve.

- To summarize the idea of Bayesian inference, If D is our data and θ is our parameter(s) Bayes' theorem tells us that

$$(2) \quad P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

- The distributions appearing in Equation (??) have the following names and interpretations:
 - $P(\theta|D)$ is the **posterior** (a beta distribution in Example ??).
 - $P(D|\theta)$ is the **likelihood** (Binomial distribution in Example ??).
 - $P(\theta)$ is the prior distribution (Uniform distribution in Example ??)
 - $P(D)$ is the evidence. It represents the chance we observe the data *unconditional* on the parameters. This can be obtained by marginalizing over the priors.

While in classical statistics, the objective is to determine (estimate) a parameter θ , in Bayesian statistics our objective is to compute the posterior distribution. Once we have this, we can obtain so-called **point estimates**, for example, by taking the average of θ or maximum of $P(\theta|D)$ (maximum likelihood).

- The beta distribution will play an important role in what follows, so we note that if

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

then (using calculus) it can be shown that

$$E[q] = \frac{a}{a+b}$$

$$\text{var}(q) = \frac{ab}{(a+b)^2(a+b+1)}.$$

Notice that this implies the estimator \hat{q}_b in Equation ?? becomes

$$\hat{q}_b = \frac{Y+1}{Y+1+N-Y+1} = \frac{Y+1}{N+2}$$

This is in-fact the estimator \hat{q}_L from the week 7 exercise set! In fact, the usual estimator \hat{q} comes from binding the value of q which maximizes $f(q|Y)$ – recall that such an estimator is called the maximum likelihood estimator.

- A very important observation, which is what the problem on the week 7 exercise set is all about, is that adding priors often has the effect of introducing a bias while reducing the variance.
- Another important feature of the β distribution is that when a and b are very large, it is approximately Normal with this mean and variance. This means we can get a rough idea of the probabilities for a β random variable. Equipped with this knowledge, we can try to go beyond the case where our prior assumption about the distribution of q , before fitting the model, is Uniform.

Example 2 (Selecting priors). *Suppose we are given a coin. We are trying to determine whether it is biased based on the outcome of N flips. We are pretty sure that, like most coins, it is not biased. To be precise, let's say you are 95% confident that the bias of the coin is less than 20% biased towards heads or tails.*

Question:

- Select a prior distribution for q .
- What is the posterior expectation if we flip the coin 5 times and see 5 heads?

Solution:

- We know that a β distribution is approximately Normal if α and β are large enough. Assuming we can make a Normal approximation, we can use the formulas for the mean and variance of the β distribution to select parameters such that

$$(3) \quad P(|q - 0.5| < 0.2) = P(0.4 < q < 0.6) \approx 0.95$$

Since we would like the mean of our prior distribution to be $1/2$, we select

$$q \sim \text{Beta}(a, a).$$

If q is approximately Normal, then Equation ?? will hold when

$$1.96\sqrt{\text{var}(q)} = 0.2$$

$$\Rightarrow \sqrt{\frac{a^2}{4a^2(2a+1)}} = \frac{1}{2}\sqrt{\frac{1}{2a+1}} = \frac{0.2}{1.6}$$

$$\Rightarrow a \approx 11.5$$

- In order to compute the posterior, we use Bayes' theorem

$$f(q|X_1, \dots, X_N) \propto q^Y (1-q)^{N-Y} \times q^{a-1} (1-q)^{a-1} = q^{Y+a-1} (1-q)^{N-Y+a-1}$$

where $Y = \sum_{i=1}^N X_i$. Again, the q dependence uniquely determines the distribution, since the area under this curve must be one. We can therefore say that

$$q|X_1, \dots, X_N \sim \text{Beta}(Y+a, N-Y+a)$$

For the mean and variance we have

$$E[q|X_1, \dots, X_N] = \frac{Y+a}{N+2a} = \frac{5+11.5}{5+2 \times 11.5} \approx 0.58$$

Note that if we used the non-Bayesian estimate we would have found $\hat{q} = 1$, while if we used uniform priors we would have found $\hat{q}_b = 0.85$.

- Something quite nice about using β -distribution priors for the Bernoulli model is that BOTH the prior and the posterior have a β -distribution, albeit with different parameters. In general, when this is the case we say that the priors are **conjugate**.

Example 3 (Bayesian inference with a Normal distribution). *Suppose we have a Normal model for Y*

$$Y|(\mu, \sigma) \sim \text{Normal}(\mu, \sigma).$$

Assume that σ is known and take our priors on μ to be

$$\mu \sim \text{Normal}(\mu_0, \sigma_\mu)$$

Question: What is the posterior distribution of μ ?

Solution: The likelihood is

$$\begin{aligned} p(D|\theta) &= \prod_i \frac{1}{\sqrt{2\sigma^2\pi}} e^{-(Y_i - \mu)^2 / 2\sigma^2} \\ &= \frac{1}{(2\sigma^2\pi)^{N/2}} e^{-\sum_i (Y_i - \mu)^2 / 2\sigma^2} \end{aligned}$$

The posterior is proportional to

$$\begin{aligned} p(D|\theta)p(\theta) &= \frac{1}{(2\sigma^2\pi)^{N/2}} e^{-\sum_i (Y_i - \mu)^2 / 2\sigma^2} \frac{1}{\sqrt{2\pi\sigma_\mu^2}} e^{-(\mu - \mu_0)^2 / 2\sigma_\mu^2} \\ &= \frac{1}{(2\sigma^2\pi)^{N/2}} \frac{1}{\sqrt{2\pi\sigma_\mu^2}} e^{-\sum_i (Y_i - \mu)^2 / 2\sigma^2 - (\mu - \mu_0)^2 / 2\sigma_\mu^2} \end{aligned}$$

On this surface this looks a bit complicated as a function of μ , but there is a trick: Notice that all the dependence on μ comes from the exponent. We can rewrite this as

$$A\mu^2 + B\mu + C$$

where

$$\begin{aligned} A &= \frac{1}{2} \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_\mu^2} \right) \\ B &= -\frac{\sum_i Y_i}{\sigma^2} - \frac{\mu_0}{\sigma_\mu^2} \\ C &= \frac{\sum_i Y_i^2}{2\sigma^2} + \frac{\mu_0^2}{\sigma_\mu^2} \end{aligned}$$

If we factor this quadratic equation, we find that it can be written

$$A(\mu - B/2A)^2 + \text{const.}$$

We don't care what the constant terms is, since it is the first term which tells us the mean and standard deviation of μ . Now observe that

$$-\frac{B}{2A} = \bar{Y} \frac{\sigma_\mu^2}{\sigma_\mu^2 + \sigma^2/N} + \mu_0 \frac{\sigma^2/N}{\sigma_\mu^2 + \sigma^2/N} \equiv \mu_b$$

and

$$\frac{1}{A} = 2 \frac{\sigma_\mu^2 \sigma^2 / N}{\sigma_\mu^2 + \sigma^2/N} \equiv 2\sigma_b^2$$

In particular, we can deduce that

$$\mu|D \sim \text{Normal}(\mu_b, \sigma_b)$$

where μ_b and σ_b^2 defined above are the mean and variance of the posterior distribution.

3. Bayesian inference for linear regression models

- We now discuss Bayesian inference for the general linear regression model with features. Let

$$f(x) = \sum_{i=1}^K \beta_i \phi_i(x)$$

and suppose that

$$\beta_i \sim \text{Normal}(0, \tau_i^2).$$

We will assume that $E[\phi_i(X)] = 0$.

- In this case, the Posterior distribution of the β_i can be understood analytically. It turns out that the marginal posterior distribution of each β_i is again Normal – thus Normal priors on β_i are conjugate priors!. The marginal mean of each β_i are determined by the system of equations given in the following Theorem.

Theorem 1 (Posterior mean for regression coefficients). Define a $K \times K$ matrix Ω called the **empirical covariance matrix** which has entries

$$\Omega_{i,j} = N \overline{\phi_i(X) \phi_j(X)} = \sum_{k=1}^N \sum_{z=1}^N \phi_i(X_k) \phi_j(X_z)$$

and let $\bar{\beta}_i$ denote the posterior mean of β_i ; that is,

$$\bar{\beta}_i = E[\beta_i|D].$$

Then $\bar{\beta}_1, \dots, \bar{\beta}_K$ satisfy the K equations

$$(4) \quad \left(\frac{\sigma}{\tau_i} \right)^2 \bar{\beta}_i + \sum_{j=1}^K \Omega_{i,j} \bar{\beta}_j = \sum_{j=1}^N \phi_i(X_j) Y_j$$

We make a few remarks on this Theorem:

- If $N > K$, then $\bar{\beta}_i$ actually have a solution. In-fact, depending on Ω , there may not be a solution. For our discussion, we will assume there is however.

- Under the assumption that $E[\phi_i(X)] = 0$, the entries of $\Omega_{i,j}$ approximate the covariances between the predictor features. Moreover, if we take the predictors to be drawn from some distribution and average over the data (both X and Y), we get

$$\frac{E[\Omega_{i,j}]}{N} = \text{cov}(\phi_j(X), \phi_i(X)).$$

and

$$\frac{1}{N} \sum_{j=1}^N E[\phi_i(X_j)Y_j] = \text{cov}(Y, \phi_i(X)).$$

Hence, we have the “math world” version of Equations ??, which is obtained by averaging over the data:

$$\frac{1}{N} \left(\frac{\sigma}{\tau_i} \right)^2 E[\bar{\beta}_i] + \sum_{j=1}^K \text{cov}(\phi_i(X), \phi_j(X)) E[\bar{\beta}_j] = \text{cov}(Y, \phi_i(X)).$$

In the limit $N \rightarrow \infty$ and with $K = 2$, we obtain a system of equations recognizable from Week 5 notes for the regression coefficients in the two predictor model in terms of the covariances. In this context, the role of the regression coefficients (which we previously took to be fixed numbers), is now played by the averages $E[\bar{\beta}_j]$, which are the average values of β_j with respect to both the posterior and the data.

- Notice that the first term in Equation ?? vanishes as either $\tau_i \rightarrow \infty$ (very broad priors) or $\sigma \rightarrow 0$ (no variance in Y conditional on the predictors). In both cases, the values of $\bar{\beta}_i$ are entirely determined by the data, and the priors play no role. When $\tau_i \rightarrow 0$ or $\sigma \rightarrow \infty$, the priors entirely determine $\bar{\beta}_i$ and in fact $\bar{\beta}_i = 0$.
- The first term in Equation ?? is an example **regularization**.
- Theorem ?? can be elegantly stated using matrices (see Linear algebra review note for intro to matrix multiplication). This is import if we want to implement these computations in Python. Observe that one way to construct $\Omega_{i,j}$ is to define the matrices

$$A = \begin{bmatrix} \phi_1(X) & \phi_2(X) & \cdots & \phi_K(X) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(X) & \phi_2(X) & \cdots & \phi_K(X) \end{bmatrix}, \quad \Lambda_0 = \begin{bmatrix} \tau_1 & 0 & \cdots & 0 \\ 0 & \tau_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \tau_K \end{bmatrix}$$

Then $\Omega = A^T A$. Now define another matrix $K \times K$ matrix

$$M = \Omega + \sigma^2 \Lambda_0^{-1}$$

Then Equations ?? can be written as

$$M\bar{\beta} = A^T y$$

where $\bar{\beta}$ and y are respectively a K -vector and N -vector

$$\bar{\beta} = \begin{bmatrix} \bar{\beta}_1 \\ \vdots \\ \bar{\beta}_K \end{bmatrix}, \quad y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix}.$$

Then $\bar{\beta}$ satisfies

$$M\bar{\beta} = A^T y \implies \bar{\beta} = (A^T A + \sigma^2 \Lambda_0^{-1})^{-1} A^T y$$