S&DS 242/542: Theory of Statistics Lecture 17: Bayesian Inference I

Bayesian inference

Given data $X_1, \ldots, X_n \stackrel{ID}{\sim} f(x \mid \theta)$ from a parametric model, we've discussed two approaches for estimating the parameter θ :

- Method-of-moments
- Maximum likelihood

We developed tools to study the sampling distributions of these estimators for large n and to construct confidence intervals for θ .

All of our discussion has treated θ as an unknown but *non-random* parameter, which describes the distribution of our data. This is called the **frequentist** paradigm of statistical inference.

In this lecture, we introduce the **Bayesian** paradigm, where θ is also modeled as random. Bayesian inference starts with our prior belief about θ , and updates this belief based on observed data.

Prior and posterior distributions

Conditional and marginal distributions

If X and Y have joint PDF or PMF $f_{X,Y}(x, y)$, then the marginal distribution of X is given by the PDF or PMF

$$f_X(x) = \int_{\mathcal{Y}} f_{X,Y}(x,y) dy, \qquad f_X(x) = \sum_{y \in \mathcal{Y}} f_{X,Y}(x,y)$$

where \mathcal{Y} is the domain of possible values for \mathcal{Y} . This describes the probability distribution of X alone, marginalizing over all possible values of Y.

The **conditional distribution** of *Y* given X = x is then defined by the PDF or PMF

$$f_{Y|X}(y \mid x) = \frac{f_{X,Y}(x,y)}{f_X(x)}$$

This describes the distribution of Y after observing X = x.

Prior distribution and likelihood

In Bayesian inference, the unknown parameter Θ is modeled as a random variable, with a probability distribution

$$\Theta \sim f_{\Theta}(heta)$$

called its **prior distribution**. This represents our prior belief about the value of Θ , before observing any data.

The parametric model or likelihood function describing the distribution of our data $\mathbf{X} = (X_1, \dots, X_n)$ is then interpreted as a conditional distribution of \mathbf{X} given Θ , which we now write as

$$\mathbf{X} \sim f_{\mathbf{X} \mid \Theta}(\mathbf{x} \mid heta)$$

(In previous lectures, we wrote this simply as $f(\mathbf{x} \mid \theta)$.)

Marginal distribution of the data

This defines a joint probability distribution over both the parameter Θ and the observed data **X**, with joint PDF/PMF

$$f_{\mathbf{X},\Theta}(\mathbf{x},\theta) = f_{\mathbf{X}|\Theta}(\mathbf{x} \mid \theta) f_{\Theta}(\theta)$$

Under this joint distribution, the marginal distribution of Θ is simply the prior $f_{\Theta}(\theta)$. The marginal distribution of the data **X** is¹

$$f_{\mathbf{X}}(\mathbf{x}) = \int f_{\mathbf{X},\Theta}(\mathbf{x},\theta) d\theta = \int f_{\mathbf{X}|\Theta}(\mathbf{x} \mid \theta) f_{\Theta}(\theta) d\theta$$

This represents the probability distribution for our observed data \mathbf{X} , if we assume Θ is also random and marginalize over its randomness.

¹If Θ is discrete, then this integral should be replaced by a sum.

Posterior distribution of the parameter

Bayesian inference is based upon the conditional distribution of the parameter Θ given the data **X**. This conditional distribution is

$$f_{\Theta \mid \mathbf{X}}(\theta \mid \mathbf{x}) = \frac{f_{\mathbf{X},\Theta}(\mathbf{x},\theta)}{f_{\mathbf{X}}(\mathbf{x})} = \frac{f_{\mathbf{X}\mid\Theta}(\mathbf{x} \mid \theta)f_{\Theta}(\theta)}{f_{\mathbf{X}}(\mathbf{x})}$$

This is called the **posterior distribution** of Θ : It represents our knowledge about the parameter Θ after observing **X**.

We often summarize the preceding equation simply as

 $f_{\Theta \mid \mathbf{X}}(\theta \mid \mathbf{x}) \propto f_{\mathbf{X} \mid \Theta}(\mathbf{x} \mid \theta) f_{\Theta}(\theta)$ Posterior \propto Likelihood \times Prior

where the symbol \propto hides the proportionality factor $1/f_{\mathbf{X}}(\mathbf{x})$ which does not depend on θ .

Example: Uniform prior for a Bernoulli proportion

Let $X_1, \ldots, X_n \stackrel{IID}{\sim}$ Bernoulli(*p*). Suppose *p* is unknown, and we model it as random with prior $P \sim \text{Uniform}(0, 1)$. To derive the posterior:

Example: Uniform prior for a Bernoulli proportion

So the mapping distribution of
$$\underline{X} = (X_{y}, X_{n})$$
 is

$$f_{\underline{X}}(X_{y}, X_{n}) = \int_{0}^{1} f_{\underline{X},p}(X_{y,p}X_{n}, p) dp$$

$$= \int_{0}^{1} p^{s}(1-p)^{n-s} dp = B(s+l, n-s+l)$$
where $B(u, \beta)$ is the Bola function

$$B(u, \beta) = \frac{\Gamma(u) \Gamma(\beta)}{\Gamma(a+\beta)} = \int_{0}^{1} p^{\alpha-1}(1-p)^{\beta-1} dp$$

$$\Rightarrow f_{\underline{P}|\underline{X}}(p|X_{u,p}X_{n}) = \frac{f_{\underline{X},p}(X_{u,p}X_{n}, p)}{f_{\underline{X}}(x)} = \frac{1}{B(s+l, n-s+l)} p^{s}(1-p)^{n-s}$$
So the perform distribution is $P \sim Bela(s+l, n-s+l), \quad s=X(r,rX_{n}, rX_{n})$

Example: Uniform prior for a Bernoulli proportion

We computed explicitly the marginal distribution $f_{\mathbf{X}}(\mathbf{x})$, but this was not actually necessary to derive the posterior distribution. Indeed, omitting factors not depending on p,

$$f_{P\mid \mathbf{X}}(p\mid \mathbf{x}) \propto f_{\mathbf{X}\mid P}(\mathbf{x}\mid p) f_{P}(p) \propto p^{s}(1-p)^{n-s}$$

This means the PDF of the posterior distribution of P has the form

$$f_{P|\mathbf{X}}(p \mid \mathbf{x}) = \frac{1}{C(\mathbf{x})} p^{s} (1-p)^{n-s}$$

where the proportionality constant $C(\mathbf{x})$ must be the unique value that makes this PDF integrate to 1 over $p \in (0, 1)$. So $C(\mathbf{x})$ is the normalizing constant for the Beta(s + 1, n - s + 1) distribution, and the posterior distribution must be Beta(s + 1, n - s + 1).

Example: Beta prior for a Bernoulli proportion

Let $X_1, \ldots, X_n \stackrel{HD}{\sim}$ Bernoulli(p). We model p as random with prior $P \sim \text{Beta}(\alpha, \beta)$. To derive the posterior distribution: Ep(p)= I por (1-p) por our p & (0,1) $f_{p|x}(p|x_{i},x_{n}) \propto f_{x|p}(x_{i},x_{n}|p) * f_{p}(p)$ $\sim p^{s}((1-p)^{n-s} * p^{s-1}((1-p)^{s-1}))$ whe S= Xititxn ~p sta-1 (1-p) n-stp-1 propertional to POP of Bota (sta, nost) > Posteriar Listibution is Pr Bota (Sta, n-Stp), S=X, t. +Xn

Example: Gamma prior in the Poisson model

Let $X_1, \ldots, X_n \stackrel{ID}{\sim} \text{Poisson}(\lambda)$. We model λ as random with prior $\Lambda \sim \text{Gamma}(\alpha, \beta)$. To derive the posterior distribution:

$$f_{\mathcal{L}}(\lambda) = \frac{\beta^{x}}{\beta(a)} \lambda^{a' \cdot i} e^{-\beta \lambda} \quad \text{for } \lambda > 0$$

$$f_{\mathcal{L}_{a}}(\lambda) = \frac{e^{-\lambda} \lambda^{\chi_{a'}}}{\chi_{a'}!}$$

$$= \int f_{A/X} (\lambda | x_{y}, x_{n}) \propto f_{X/A} (x_{y}, x_{n} | \lambda) \times f_{A}(\lambda)$$

$$= \left(\prod_{i=1}^{n} e^{-\lambda} \lambda^{x_{i}} \right) \times \lambda^{\alpha-i} e^{-\beta \lambda}$$

$$= \left(\sum_{i=1}^{n} e^{-(\alpha+\beta)\lambda} - \sum_{i=1}^{n} e^{-(\alpha+\beta)\lambda} \right)$$

$$= \left(\sum_{i=1}^{n} e^{-(\alpha+\beta)\lambda} - \sum_{i=1}^{n} e^{-\lambda} \right)$$

$$= \int f_{A/X} \left(\sum_{i=1}^{n} e^{-\lambda} \right) + \int f_{A/X} \left$$

Conjugate priors

A family of priors $f_{\Theta}(\theta)$ are called **conjugate priors** for a parametric model $f_{\mathbf{X}|\Theta}(\mathbf{x} \mid \theta)$ if the resulting posteriors $f_{\Theta|\mathbf{X}}(\theta \mid \mathbf{x})$ have the same parametric form as the priors. So

• Beta (α, β) are conjugate priors for the model Bernoulli(p)

• Gamma(α, β) are conjugate priors for the model Poisson(λ) Conjugate priors exist for many simple parametric models, and are mathematically/computationally convenient because they lead to posteriors that are common distributions and easy to compute.

This convenience does not necessarily imply desirable statistical properties — for example, conjugate priors are typically light-tailed distributions that bias posterior inferences towards the prior mean. Non-conjugate (in particular heavier-tailed) priors are also often used to yield posterior inferences that are more robust to the choice of the prior, even if this requires more intensive computation.

Example: Normal prior for a normal mean

Let $X_1, \ldots, X_n \stackrel{ID}{\sim} \mathcal{N}(\theta, \frac{1}{\xi})$, where for convenience we parametrize the model by the *precision* $\xi = \frac{1}{\sigma^2}$ instead of the variance σ^2 . Suppose ξ is known, θ is unknown with prior $\Theta \sim \mathcal{N}(\mu_{\text{prior}}, \frac{1}{\xi_{\text{prior}}})$.

$$\begin{aligned} f_{(i)}(\theta) &= \sqrt{\frac{5}{2\pi}} e^{-\frac{5}{2\pi}} (\theta - M_{i} L_{i})^{2} \\ f_{X_{i}}(\theta) &= \sqrt{\frac{5}{2\pi}} e^{-\frac{5}{2}} (x_{i} - \theta)^{2} \\ \Rightarrow f_{\theta|X}(\theta|x_{i}, x_{i}) \propto f_{X|\theta}(x_{i}, x_{i}|\theta) * f_{\theta}(\theta) \\ \propto (\frac{1}{11} e^{-\frac{5}{2}} (x_{i} - \theta)^{2})^{2} e^{-\frac{5}{2\pi}} (\theta - M_{i}L_{i})^{2} \\ \propto e_{y_{0}} \left(-\frac{5}{2\pi} \frac{5}{2} (x_{i} - \theta)^{2} - \frac{5}{2\pi} \frac{1}{2} (\theta - M_{i}L_{i})^{2}\right) \end{aligned}$$

Example: Normal prior for a normal mean

× exp (- 2 (+ 0 - 5 × 0) - (+ 0 - 5 × 0)) acpp (- 2 (n 3+ 5, 1) 02+ (5 2 × + 5, 1) 0)

Lot Sport = n 3+ 3pin Mp. 1 = g = x + Spin Mpin =) folx (O(K, x) & cyp (- 50.4 02 + 50.50 More O) x cxp (- 52 (0-1)2) => Poster'ar is @~N(Mpost, \$post)

Example: Normal prior for a normal variance

Let $X_1, \ldots, X_n \stackrel{ID}{\sim} \mathcal{N}(\theta, \frac{1}{\xi})$. Suppose θ is known, ξ is unknown with prior $\Xi \sim \text{Gamma}(\alpha, \beta)$.

$$f_{\underline{r}}(\underline{s}) = f_{\underline{r}(\underline{s})} \underline{s}^{\alpha''} e^{-\beta \underline{s}}, \quad \underline{s} > 0$$

$$f_{\underline{x}_1|\underline{s}}(\underline{x}_1|\underline{s}) = \sqrt{\underline{z}_{\overline{n}}} e^{-\underline{z}(\underline{x}_1 - \theta)^{\mathsf{T}}}$$

$$\Rightarrow f_{\underline{s}|\underline{x}}(\underline{s}|\underline{x}_1,\underline{s}) \propto f_{\underline{x}|\underline{s}}(\underline{x}_1,\underline{s},\underline{x}_1|\underline{s}) \times f_{\underline{s}|\underline{s}}(\underline{s})$$

$$\propto (f_{\underline{r}|1}^{\mathsf{T}} \int \underline{s} e^{-\frac{\underline{s}}{\underline{s}}(\underline{x}_1 - \theta)^{\mathsf{T}}}) \times \underline{s}^{\alpha''} e^{-\beta \underline{s}}$$

$$\propto \underline{s}^{\alpha' + \frac{\eta}{2} - 1} e^{-\frac{\underline{s}}{\underline{s}}(\underline{x}_1 - \theta)^{\mathsf{T}}} e^{-\beta \underline{s}}$$

$$\Rightarrow Posture i_{s} = \frac{1}{s} \wedge Gamm_{s}(\alpha' + \frac{\eta}{2}, \frac{1}{z}, \frac{z}{z})(\underline{x}_1 - \theta)^{\mathsf{T}} + \beta)$$

Point estimates and credible intervals

Bayesian point estimates and credible intervals

In many applications, we may want a single estimate $\hat{\theta}$ and an interval that quantifies our uncertainty about θ , instead of a full posterior distribution.

The **posterior mean** and **posterior mode** are the mean and mode of the posterior distribution of Θ , and either may be used as a Bayesian estimate $\hat{\theta}$ for θ .

A **Bayesian credible interval** with coverage $1 - \alpha$ is an interval $I(\mathbf{X})$ that contains Θ with posterior probability $1 - \alpha$:

$$\mathbb{P}[\Theta \in I(\mathbf{X}) \mid \mathbf{X} = \mathbf{x}] = 1 - \alpha$$

where the probability is over Θ , instead of **X** as in the frequentist confidence interval. A common choice is the interval from the lower- $\alpha/2$ to upper- $\alpha/2$ point of the posterior distribution of Θ .

Example: Beta prior for a Bernoulli proportion

Let $X_1, \ldots, X_n \stackrel{HD}{\sim}$ Bernoulli(*p*), with prior $P \sim \text{Beta}(\alpha, \beta)$. Recall that the posterior distribution is

$$P \sim \text{Beta}(S + \alpha, n - S + \beta), \qquad S = X_1 + \ldots + X_n$$

The posterior mean is

$$\hat{p} = \frac{S + \alpha}{n + \alpha + \beta}$$

This is different from the MLE/MoM-estimator $\bar{X} = \frac{S}{n}$, and is biased towards the prior mean.

One interpretation is that \hat{p} is the sample mean as if we had observed — a priori — α heads and β tails before seeing our data. This gives an interpretation also for the Beta(α, β) prior.

Example: Beta prior for a Bernoulli proportion

An alternative interpretation is that \hat{p} is a weighted average of the sample mean and prior mean,

$$\hat{p} = \frac{S + \alpha}{n + \alpha + \beta} = \underbrace{\frac{n}{n + \alpha + \beta}}_{\text{sample weight sample mean}} \cdot \underbrace{\frac{S}{n}}_{\text{sample mean}} + \underbrace{\frac{\alpha + \beta}{n + \alpha + \beta}}_{\text{prior weight prior mean}} \cdot \underbrace{\frac{\alpha}{\alpha + \beta}}_{\text{prior mean}}$$

For any fixed prior $Beta(\alpha, \beta)$, as *n* increases, the sample weight increases to 1 and prior weight decreases to 0.

The interval from the lower-0.05 point to the upper-0.05 point of the Beta($S + \alpha$, $n - S + \beta$) posterior distribution forms a 90% Bayesian credible interval for p.