

BAYES ESTIMATION

This handout is adapted from one for 18.443 and a section from previous 18.466 notes (2005 version, a revision of the 18.466 OCW 2003 notes).

1. DEFINITION OF PRIORS AND POSTERIOR FOR A CONTINUOUS θ

In this handout Θ will be a parameter space included in a Euclidean space \mathbb{R}^k . For example, for the family of normal distributions, Θ is the open half-plane $\{(\mu, \sigma) : -\infty < \mu < \infty, 0 < \sigma < \infty\} \subset \mathbb{R}^2$. On Θ , $d\theta$ will mean $d\theta_1 \cdots d\theta_k$.

Assume given a likelihood function $f(X, \theta)$ defined for $\theta \in \Theta$ and X a vector in \mathbb{R}^n . In Bayesian statistics, one assumes before taking any observations that θ has a prior probability density $\pi(\theta)$ with respect to $d\theta$. Sometimes, as in some examples in Bickel and Doksum, “improper” priors with infinite total measure such as Lebesgue measure on the whole line are considered. This handout treats only prior probability measures.

For an ordinary (proper) prior, $\pi(\theta) \geq 0$ and $\int_{\Theta} \pi(\theta) d\theta = 1$. If $\theta = p$ with $0 \leq p \leq 1$ is the success probability in a binomial distribution, a simple and natural choice for its prior (in the absence of any particular information about p) is a $U[0, 1]$ distribution with $\pi(p) = 1$ for $0 \leq p \leq 1$. The earliest works in Bayesian statistics, Bayes (1764) and Laplace (1774), made this choice.

Let $f(x, \theta)$ be a likelihood function for one observation, which may be either a probability mass function if x is discrete or a density function if x is continuous. If we have i.i.d. observations $X = (X_1, \dots, X_n)$ we get a likelihood function $f(X, \theta) = \prod_{j=1}^n f(X_j, \theta)$.

However $f(X, \theta)$ is obtained, the *posterior* density $\pi_X(\theta)$ is gotten by multiplying the likelihood function by the prior and then normalizing it,

$$(1) \quad \pi_X(\theta) = \frac{f(X, \theta)\pi(\theta)}{\int_{\Theta} f(X, \phi)\pi(\phi)d\phi}.$$

To show that (1) makes sense we can use the following (in it, to be measure-theoretically accurate, it should be assumed that f is a jointly measurable function of X and θ):

Theorem 1. *Let $\Theta \subset \mathbb{R}^k$ be a parameter space, $\theta \in \Theta$, and $X \in \mathbb{R}^n$ an observed vector. Suppose that for each $\theta \in \Theta$, $f(X, \theta)$ is a probability density with respect to X , so that $\int f(X, \theta) dX = 1$ where $dX = dx_1 dx_2 \cdots dx_n$. Let $\pi(\theta) \geq 0$ be a prior probability density for θ . Let $q(X, \theta) = \pi(\theta)f(X, \theta)$ for all $\theta \in \Theta$ and all X . Then*

(a) *q is a bivariate probability density with respect to $d\theta dX$, for a joint probability distribution Q of (X, θ) ,*

(b) *the marginal density of q with respect to θ is π ,*

(c) *and for each $\theta \in \Theta$ the conditional density of X given θ is $q(X|\theta) = f(X, \theta)$.*

(d) *Letting*

$$\tau(X) = \int_{\Theta} q(X, \theta) d\theta,$$

τ is a probability density and is the marginal density of Q with respect to X .

(e) *With probability 1 with respect to Q , or with respect to its marginal density τ ,*

$$(2) \quad 0 < \tau(X) < +\infty.$$

(f) For all X such that (2) holds, a conditional density of θ given X exists and is given by $q(\theta|X) = \pi_X(\theta)$ in (1) where the denominator in (1) is $\tau(X)$.

(g) We have for Q -almost all (x, θ) ,

$$(3) \quad q(X, \theta) = \pi(\theta)f(X, \theta) = \tau(X)\pi_X(\theta).$$

Remark. The theorem adapts easily to the case where X is discrete, so $f(\cdot, \theta)$ is a probability mass function, with integrals $\int \cdot dX$ replaced by sums \sum_X .

Proof. For (a), since the integrand is nonnegative we can do an iterated integral in either order. If we integrate first with respect to X we get $\pi(\theta)$ which has integral 1 with respect to θ . This also proves (b), and the rest of the statements are known facts about marginal and conditional densities from probability theory. Since part (e) is crucial in showing that π_X is well-defined with probability 1, let's prove it in detail, assuming part (d). We have

$$\Pr(\tau(X) = 0) = \int_{\tau(X)=0} \tau(X) dX = \int 0 dX = 0.$$

On the other hand let $A = \{X : \tau(X) = +\infty\}$. Then

$$\Pr(A) = \int_A \tau(X) dX = \int_A +\infty dX = \Pr(A) \cdot (+\infty) = +\infty$$

if $\Pr(A) > 0$, but since $\Pr(A) \leq 1$, we get $\Pr(A) = 0$, and (e) follows, i.e. (2) holds with probability 1.

For part (g), in (3), the first equation holds by definition of $q(X, \theta)$, and the second by the definitions and parts (d) and (e). \square

2. CONJUGATE PRIORS

A conjugate prior for a given parametric family of distributions with a likelihood function is one such that the posterior distributions all belong to the same parametric family. For example, if $\theta = \lambda$ is a Poisson parameter with $0 < \lambda < +\infty$ and the prior $\pi(\theta)$ is a gamma density, then the posterior $\pi_X(\theta)$ is also in the gamma family. Specifically, if λ has prior density $\text{Gamma}(a, c)$, where $a > 0$ and $c > 0$, so that $\pi(\lambda) = c^a \lambda^{a-1} \exp(-c\lambda) / \Gamma(a)$ and we observe X_1, \dots, X_n i.i.d. $\text{Poisson}(\lambda)$ with $S_n := X_1 + \dots + X_n$, then the likelihood function is proportional to $e^{-n\lambda} \lambda^{S_n}$ and so the posterior density is $\Gamma(a + S_n, c + n)$ (it is proportional to this as a function of λ , and a probability density has a unique normalizing constant). As the expectation for $\Gamma(a, c)$ is a/c , the expectation for the posterior distribution is $\frac{S_n + a}{n + c}$, which is asymptotic as $n \rightarrow \infty$ to the maximum likelihood estimate S_n/n .

Likewise, beta densities give conjugate priors for the binomial probability p .

2.1. Normal-inverse-gamma distributions; conjugate for normals. Let $Y > 0$ be a random variable having a distribution function F and a density $f = f_Y$. Let $V := 1/Y$. Then for any $x > 0$, $\Pr(V \leq x) = \Pr(1/Y \leq x) = \Pr(Y \geq 1/x) = 1 - F(1/x)$, and so

by the chain rule $1/Y$ has a density $f_{1/Y}(x) = -f(1/x) \cdot (-1/x^2) = x^{-2}f(1/x)$. Thus if Y has a Gamma(α, β) density $f(y) = \beta^\alpha y^{\alpha-1} \exp(-\beta y)/\Gamma(\alpha)$ then $1/Y$ has the density

$$\beta^\alpha y^{1-\alpha} \exp(-\beta/y)/(y^2\Gamma(\alpha)) = \beta^\alpha y^{-1-\alpha} \exp(-\beta/y)/\Gamma(\alpha)$$

where $\alpha > 0$, $\beta > 0$, and $y > 0$, which is called an inverse gamma(α, β) density.

Parameters of prior or posterior distributions are called hyperparameters. The family of all normal distributions $N(\mu, \sigma^2)$ on the real line has a conjugate prior for the parameter $\theta = (\mu, \sigma^2)$ called the “normal-inverse-gamma distribution” and given by

$$(4) \quad \frac{\sqrt{\nu}}{\sigma\sqrt{2\pi}} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left(-\frac{2\beta + \nu(\mu - \lambda)^2}{2\sigma^2}\right)$$

with four hyperparameters $\alpha > 0$, $\beta > 0$, $\nu > 0$, and λ which can be any real number. For (4), given the hyperparameters, the marginal density of σ^2 is inverse gamma (α, β), or equivalently $1/\sigma^2$ has Gamma(α, β), and the conditional density of μ given σ is $N(\lambda, \sigma^2/\nu)$. If σ is fixed, then the normal distributions give a conjugate prior family for μ , which is much simpler, but it's usually unrealistic to assume σ is known. Likewise if μ is fixed, the gamma distributions for $1/\sigma^2$ give a conjugate prior family, but for μ to be fixed is also usually unrealistic. For the joint conjugate prior density (4) of μ and σ^2 , μ and σ^2 are not independent: the density is not a product $f(\mu)$ times $g(\sigma)$ for any functions f and g . So the joint conjugate prior is a bit complicated.

3. CREDIBLE INTERVALS

These are the Bayesian counterparts of confidence intervals. A $100(1 - \alpha)\%$ credible interval for a real parameter θ is one that has posterior probability $1 - \alpha$ of containing θ . A two-sided 95% credible interval for θ , for example, would be the interval with endpoints the 0.025 and 0.975 quantiles of the posterior distribution.

4. BAYES LEAST-SQUARES ESTIMATION

First here is a very simple fact.

Proposition 2. *For any random variable Y with $E(Y^2) < +\infty$, the unique constant c that minimizes $E((Y - c)^2)$ is $c = EY$.*

Proof. $E((Y - c)^2) = E(Y^2) - 2cEY + c^2$ is a quadratic polynomial in c which goes to $+\infty$ as $c \rightarrow \pm\infty$, so it's minimized where its derivative with respect to c equals 0, namely at $c = EY$, Q.E.D.

Suppose we want to estimate a function $g(\theta)$. Then for an estimator $V(X)$, the mean-square error (MSE) for a given θ is $E_\theta[(V(X) - g(\theta))^2]$. For a prior π , the *risk* is the expectation of the MSE with respect to that prior, namely

$$(5) \quad r(V, \pi) := \int_{\Theta} E_\theta[(V(X) - g(\theta))^2] \pi(\theta) d\theta.$$

A *Bayes* estimator for $g(\theta)$ for the given prior is one that minimizes the risk, provided its risk is finite.

Theorem 3. For a given likelihood function $f(X, \theta)$ for $\theta \in \Theta$, where $\Theta \subset \mathbb{R}^k$ for some $k \geq 1$, and prior density π , if there exists some estimator $U(X)$ of the given $g(\theta)$ that has finite risk for the given π , then there exists a Bayes estimator T , given by the expectation of $g(\theta)$ with respect to the posterior distribution,

$$(6) \quad T(X) = \int_{\Theta} g(\theta) \pi_X(\theta) d\theta.$$

The Bayes estimator is essentially unique, in the sense that any Bayes estimator must equal this $T(X)$ with probability 1.

Proof. We would like to minimize (5). Let's write out the E_{θ} . Recall that $\int \cdots dX$ is a shorthand for

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \cdots dx_1 dx_2 \cdots dx_n,$$

where the integral(s) are replaced by sums in case X is discrete. The method of proof is essentially the same. Then (5) becomes

$$(7) \quad \int_{\Theta} \int [(V(X) - g(\theta))^2] f(X, \theta) dX \pi(\theta) d\theta.$$

Since the integrand is nonnegative and the integrals are well-defined (possibly infinite) we can interchange the two integrals, and (7) becomes

$$(8) \quad \int \int_{\Theta} [(V(X) - g(\theta))^2] f(X, \theta) \pi(\theta) d\theta dX.$$

Then applying (3), the factor $\tau(X)$ doesn't depend on θ so we can take it outside the integral with respect to θ , and (8) becomes

$$(9) \quad \int \int_{\Theta} [(V(X) - g(\theta))^2] \pi_X(\theta) d\theta \tau(X) dX.$$

In the inner integral with respect to θ in (9), X is fixed and $g(\theta)$ is a random variable with respect to the posterior density $\pi_X(\theta)$. To minimize this inner integral we need to choose $V(X)$, which would be constant for fixed X . By Proposition 2, the correct constant is given by $V(X) = T(X)$ in (6). Since the risk is finite for some estimator by assumption, the minimum risk must be finite, so $T(X)$ in (6) indeed gives a Bayes estimator. The essential uniqueness follows from the uniqueness in Proposition 2. Q.E.D.

In case of a Gamma(a, c) prior density for a Poisson parameter λ , where the posterior density will also be in the gamma family, the expectation of λ for the posterior density is easy to calculate, as we saw above. Similarly, we have an easy calculation for the posterior expectation of a binomial parameter p using a Beta(a, b) prior.

Some texts give a different formulation of Theorem 3 in which they say that the Bayes estimator is the conditional expectation of $g(\theta)$ given X , $T(X) = E(g(\theta)|X)$. That is correct in case $\int |g(\theta)| \pi(\theta) d\theta < +\infty$, but integrals with respect to the posterior distributions may be finite even if they are not with respect to the prior, as will be seen in a problem. There is more about conditional expectations in the Notes at the end.

5. ADMISSIBILITY

Recall that a statistic $T(X)$ is said to be *inadmissible* as an estimator of a function $g(\theta)$ of a parameter θ if there exists another estimator $V(X)$ such that $E_\theta((V(X) - g(\theta))^2) \leq E_\theta((T(X) - g(\theta))^2)$ for all θ and $E_\theta((V(X) - g(\theta))^2) < E_\theta((T(X) - g(\theta))^2)$ for some θ . Then $T(X)$ is *admissible* if it is not inadmissible. Let's call $T(X)$ *strongly inadmissible* if we add to the definition that $E_\theta[(V(X) - g(\theta))^2] < E_\theta[(T(X) - g(\theta))^2]$ for all θ in a non-empty open set U , namely, a set such that: for some θ_0 in U and $r > 0$, also θ is in U for all θ such that $|\theta - \theta_0| < r$. In one dimension this would just say that U includes a non-degenerate interval.

If π is a prior density with $\pi(\theta) > 0$ for almost all θ , i.e. if A is the set of θ for which $\pi(\theta) > 0$, then $\int 1_A(\theta)d\theta = 0$, and if T is a Bayes estimator for $g(\theta)$, namely the integral of $g(\theta)$ times the posterior density $\pi_X(\theta)$, then T cannot be strongly inadmissible, or there would be an estimator with smaller overall risk (integrating mean-square error times $\pi(\theta)$), contradicting the Bayes property of T .

6. UNBIASEDNESS

The Bayes property for squared-error loss turns out to be virtually incompatible with unbiasedness. Let's begin with two examples.

Example 1. For $0 \leq \theta \leq 1$ let δ_θ be the point mass at θ . Suppose θ is unknown in advance and has prior density $U[0, 1]$. Suppose that the true $\theta = \theta_0$ for some θ_0 . If we have even one observation X_1 from δ_θ , then $\Pr(X_1 = \theta_0) = 1$. So X_1 is both an unbiased estimator and a Bayes estimator of θ for the given prior (or any prior on $[0, 1]$). That shows that this combination of properties of an estimator is possible, but it's a rather extreme and impractical case.

Example 2. Let p be the success probability in a binomial(n, p) distribution, and let $\pi(p) > 0$ for $0 < p < 1$ be a prior density for p . Then for any observation X , an integer with $0 \leq X \leq n$, the likelihood function is proportional to $p^X(1-p)^{n-X}$, which is a bounded function of p , and a posterior density $\pi_X(p)$ defined by (1) exists. Moreover, since p is bounded, the integral $T(X) = \int_0^1 p\pi_X(p)dp$ is always finite, in fact satisfies $0 \leq T(X) \leq 1$, and so has finite risk as an estimator of p with squared-error loss. So by Theorem 3, T is a Bayes estimator for p for the given π . Now suppose the true $p = 0$. Then we will have $\Pr(X = 0) = 1$ and the likelihood function will be $(1-p)^n$. Then $\pi_X(p) > 0$ for $0 < p < 1$ because $\pi(p) > 0$, and so $\Pr(T(X) > 0) = 1$ and $E_0T(X) > 0$, so T is not an unbiased estimator of p .

Similarly, when the true $p = 1$, $E_1(T(X))$ will be less than 1. For $0 < p_0 < 1$ there do exist priors π of p such that for the Bayes estimator $T(X)$ for π , $E_{p_0}T(X) = p_0$. But it is not possible to find π such that for the Bayes estimator $T(X)$ for π , $E_p(T(X)) = p$ for all p with $0 < p < 1$, as a special case of the following theorem.

Theorem 4. *Let $f(X, \theta)$, $\theta \in \Theta$, be a parametric family of densities for X in n -dimensional Euclidean space \mathbb{R}^n , with respect to $dX = dx_1dx_2 \cdots dx_n$, where θ is in any parameter space Θ included in a Euclidean space \mathbb{R}^k . For any prior density π on Θ and real-valued function g on Θ which is a random variable with respect to π , an*

unbiased estimator T of g is Bayes for π and squared-error loss if and only if it has risk $r(T, \pi) = 0$, so that $T(x) = g(\theta)$ with Q -probability 1.

Remark. The theorem shows that an estimator can be both Bayes and unbiased only when $g(\theta)$ can be estimated exactly without error, as in Example 1.

Proof. “If” is clear. To prove “only if,” by definition of Bayes estimator, T must have finite risk. Let τ be the marginal density of Q for X given by Theorem 1(d). By (3), $q(X, \theta) = \pi_X(\theta)\tau(X)$ with probability 1. We have

$$\begin{aligned} r(T, \pi) &= \int \int (T(X) - g(\theta))^2 q(X, \theta) dX d\theta \\ (10) \quad &= \int \int T(X)^2 - 2T(X)g(\theta) + g(\theta)^2 \pi_X(\theta)\tau(X) dX d\theta. \end{aligned}$$

As the integrand in (10) is nonnegative we can do the integral in either order. The proof will work by finding two different expressions of the integral of the cross term $-2T(X)g(\theta)$. For one of them, by the Bayes property and equation (6),

$$T(X) = \int g(\theta)\pi_X(\theta)d\theta.$$

Doing the integral in (10) in the order $d\theta dX$ and doing the integral with respect to θ of the cross term, X is fixed and we get $-2T(X)\tau(X)T(X) = -2T(X)^2\tau(X)$. It then follows that

$$(11) \quad r(T, \pi) = \int \left[T(X)^2 - 2T(X)^2 + \int g(\theta)^2 \pi_X(\theta) d\theta \right] \tau(X) dX.$$

Since $r(T, \pi) < \infty$, and for fixed X , $-T(X)^2$ is also fixed, we have $\int g(\theta)^2 \pi_X(\theta) d\theta < \infty$ for τ -almost all X , and

$$(12) \quad r(T, \pi) = \int \int [g(\theta)^2 - T(X)^2] q(X, \theta) d\theta dX.$$

On the other hand, doing the integral in (10) in the stated order, we know by unbiasedness that for fixed θ , $E_\theta T(X) = \int T(X)f(X, \theta)dX = g(\theta)$. By (3) $f(X, \theta)\pi(\theta) \equiv \pi_X(\theta)\tau(X)$. As θ is fixed in the inner integral dX , we can take $\pi(\theta)$ outside the integral. We then have

$$\int -2g(\theta)T(X)f(X, \theta)dX = -2g(\theta)^2$$

and so

$$r(T, \pi) = \int \left[\int T(X)^2 f(X, \theta) dX - 2g(\theta)^2 + g(\theta)^2 \right] \pi(\theta) d\theta.$$

Next, $r(T, \pi) < \infty$ implies $\int T(X)^2 f(X, \theta) dX < \infty$ for π -almost all θ , and

$$r(T, \pi) = \int [T(X)^2 - g(\theta)^2] q(X, \theta) dX d\theta = -r(T, \pi)$$

from (12), so $r(T, \pi) = 0$, finishing the proof. \square

There are cases where a maximum likelihood estimate (MLE) is unbiased, as with the sample mean \bar{X} for the normal mean μ or the Poisson parameter λ . In such cases, typically a Bayes estimator will be somewhere between the MLE and the mean of the prior distribution, becoming asymptotic to the MLE as $n \rightarrow \infty$.

7. NOTES

7.1. Conditional expectations. In measure and probability theory, usually conditional expectation is taken only of functions X which have a finite expectation. For example, a martingale is defined as a sequence $\{X_n\}_{n \geq 0}$ of random variables with finite expectations EX_n (where $E|X_n| < +\infty$ for all n) together with an increasing sequence of σ -algebras $\{\mathcal{F}_n\}_{n \geq 0}$ such that each X_j is measurable with respect to \mathcal{F}_j and the conditional expectation $E(X_{n+1}|\mathcal{F}_n) = X_n$ for all n , e.g. Dudley (2002, Chapter 10). But, some statisticians sometimes write conditional expectations $E(g|T)$ given a statistic T where g may not have a finite (unconditional) expectation, as we'll see.

Theorems 3 and 4 are stated in Lehmann (1991), Corollary 4.1.1 p. 239 and Theorem 4.1.2 pp. 244–245, but for the former, Lehmann writes the Bayes estimator as $T(x) = E(g(\theta)|x)$. Lehmann's proof of the latter theorem uses the assumption that $g(\theta) \in \mathcal{L}^2(\pi)$, in other words $\int g(\theta)^2 d\pi(\theta) < +\infty$. We will see that Bayes estimators may exist when even $\int |g(\theta)| d\pi(\theta)$ may be infinite. Thus, other proofs have been given for Theorems 3 and 4 without any moment assumptions.

Lehmann apparently doesn't give any earlier references for these facts, although at least Theorem 3 for $g \in \mathcal{L}^2$ was presumably known well before 1983.

Bickel and Doksum, Second ed., (3.2.5) p. 162, assert that either all estimators of $g(\theta)$ have infinite risk for squared-error loss, or if there is one with finite risk, then the Bayes estimator of $g(\theta)$ is $E(g(\theta)|x)$.

One might say perhaps that some statisticians' definition of $E(g(\theta)|x)$ is the expectation of $g(\theta)$ with respect to the conditional *distribution* of θ given x , which exists and is the posterior distribution.

Berger (1985, Section 4.4.2 p. 161) states that the Bayes estimator for squared-error loss is the expectation for the posterior distribution, in the special case $g(\theta) \equiv \theta$, under the assumption that each of three integrals for the posterior distribution is finite (as they will be, almost surely, under the assumption of Theorem 3). Such a statement avoids any ambiguity about conditional expectations.

7.2. Historical Notes (18th century origins). These notes are based on Stigler (1986), pp. 359–362. The field of “Bayesian” statistics is named for Thomas Bayes, who wrote a paper about the method in 1764. But the work by the leading mathematician and scientist Laplace (1774) attracted more attention. Stigler (p. 361) wrote that Bayes's article “was ignored until after 1780 and played no important role in scientific debate until the twentieth century.” Laplace used the $U[0, 1]$ prior distribution for a binomial parameter p and noted that the posterior distributions are beta distributions. Stigler (p. 359) writes “we can be reasonably certain Laplace was unaware of Bayes's earlier work.”

REFERENCES

- Bayes, T. (1764), An essay toward solving a problem in the doctrine of chances, *Philos. Trans. Roy. Soc. London* **53**, 370-418; repr. in *Biometrika* **45** (1958) 293-315.
- Berger, J. O. (1985), *Statistical Decision Theory and Bayesian Analysis*, Springer, New York. (Second Ed. of *Statistical Decision Theory*, 1980, Springer.)
- Dudley, R. M. (2002). *Real Analysis and Probability*, second ed. Cambridge University Press, New York.
- Dudley, R. M. (2003). *Mathematical Statistics*, 18.466 lecture notes, Spring 2003. On MIT OCW (OpenCourseWare) website, 2004.
- Laplace, P. S. (1774), "Memoir on the probability of the causes of events," *Mémoires de Math. et de Physique, Présentés à l'Académie Royale des Sciences, par divers Savans & lus dans ses Assemblées*, **6**, 621-646; transl. by S. M. Stigler in Stigler (1986), pp. 364-378.
- Stigler, S. M. (1986), Laplace's 1774 memoir on inverse probability, *Statistical Science* **1**, 359-378.