

Some Asymptotic Bayesian Inference

(background to Chapter 2 of Tanner's book)

Principal Topics

- The approach to *certainty* with increasing evidence.
- The approach to *consensus* for several agents, with increasing shared evidence.
- A role for *statistical models* in these asymptotic results.
 - symmetry/independence assumptions in these results.
 - data reduction
 - asymptotic Normal inference for these results.

Generalizing the *coin-tossing example* from last lecture:

Sample space of (observable) outcomes:

A 2-sided coin is repeatedly tossed, indefinitely,

$$X = \langle X_1, X_2, \dots, X_n, \dots \rangle$$

$X_j = 0$, or $X_j = 1$ as the coin lands *tails up* or *heads up* on the j^{th} flip.

So that, $\mathbf{x} = \langle x_1, x_2, \dots, x_n, \dots \rangle$ is a point of the space $\Omega = \{0,1\}^{\aleph_0}$

Of course, at any one time we observe only a finite, initial segment.

The *events* that make up the σ -algebra, \mathbf{A} , are the (smallest) σ -field of sets including all the (historical) observable events, of the form,

$$\mathbf{H}_n = \langle x_1, x_2, \dots, x_n, \{0,1\}, \{0,1\}, \dots \rangle$$

The Statistical Model:

Introduce a statistical quantity, a parameter θ , such that the events in \mathbf{A} have a determinate conditional probability, given the parameter.

Bernoulli (i.i.d.) Coin flipping (continued):

$$\mathbf{P}(X_j = 1 | \theta) = \theta \quad (j = 1, \dots), \text{ for } 0 \leq \theta \leq 1$$

$$\mathbf{P}(\mathbf{H}_n | \theta) = \theta^k (1-\theta)^{n-k}, \text{ where } k \text{ of the first } n \text{ coordinates of } \mathbf{H}_n \text{ are } 1$$

and $n-k$ of the first n coordinates \mathbf{H}_n of are 0.

Now, if we are willing to make θ into a random variable (by expanding the σ -algebra accordingly), we can write Bayes theorem for the parameter:

$$\begin{aligned} \mathbf{P}(\theta | H_n) &= \mathbf{P}(H_n | \theta) \mathbf{P}(\theta) / \mathbf{P}(H_n) \\ &\propto \mathbf{P}(H_n | \theta) \mathbf{P}(\theta) \end{aligned}$$

OR

The *posterior probability for θ* is proportional to

the product of the *likelihood for θ* and its *prior probability*.

With the conjugate **Beta**(α, β) prior for θ

$$\mathbf{P}(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\alpha, \beta > 0, 0 \leq \theta \leq 1,$$

the posterior distribution $\mathbf{P}(\theta | \mathbf{H}_n)$ is given by the distribution **Beta**($\alpha+k, \beta+n-k$)

having *mean* $(\alpha+k) / (\alpha+k+\beta+n-k) = (\alpha+k) / (\alpha+\beta+n)$

and *variance* $(\alpha+k)(\beta+n-k) / (\alpha+\beta+n)^2(\alpha+\beta+n+1)$

Note: Here we may reduce the historical data of n -bits to two quantities ($k, n-k$). That is, the two **likelihood functions**: $\mathbf{P}(\mathbf{H}_n | \theta)$ and $\mathbf{P}(k, n-k | \theta)$ are the same.

$$\mathbf{P}(\mathbf{H}_n | \theta) = \mathbf{P}(k, n-k | \theta)$$

Now, by the *Strong Law of Large Numbers*: for each $\varepsilon > 0$ and given θ

$$\mathbf{P}(\lim_{n \rightarrow \infty} |k/n - \theta| < \varepsilon \mid \theta) = 1.$$

Hence, with probability 1, the sequence of posterior probabilities for θ

$$\lim_{n \rightarrow \infty} \mathbf{P}(\theta \mid \mathbf{H}_n) = \lim_{n \rightarrow \infty} \mathbf{Beta}(\alpha + n\theta, \beta + n(1 - \theta))$$

have a limit distribution with *mean* θ and *variance* 0, independent of α and β .

Note: The posterior variances for θ are $O(1/n)$. That is, in advance, we can bound from below the precision (that is, bound from above the variance) of the posterior distribution for the parameter by choosing the sample size to observe.

Asymptotic Certainty

THUS, under each conjugate (Beta) prior:

With probability 1, the posterior probability for θ converges to the (0-1) *Delta* distribution, concentrated on the *true* parameter value.

Aside: (Doob, 1949) if the parameter space is *finite-dimensional*, this almost sure convergence occurs for **each** value of the parameter in the *support* of the prior.

(The *support* of the prior in the parameter set is the smallest closed set with prior prob. 1)

From the perspective of the posterior probability for θ ,
through the likelihood function, the data “swamp” the prior.

Asymptotic Consensus (Merging of Posterior Probability Distributions)

As a metric (distance) between two distributions \mathbf{P} and \mathbf{Q} over the algebra A , consider a strict standard, *uniform distance*,

$$\rho(\mathbf{P}, \mathbf{Q}) = \sup_{E \in A} | \mathbf{P}(E) - \mathbf{Q}(E) |$$

Let $\mathbf{P}^n = \mathbf{P}(\theta | \mathbf{H}_n)$ and $\mathbf{Q}^n = \mathbf{Q}(\theta | \mathbf{H}_n)$ ($n = 1, 2, \dots$) be two sequences of posterior probability distributions for the parameter θ based on two (conjugate) Beta priors.

Then, it is not hard to show that

$$\lim_{n \rightarrow \infty} \rho(\mathbf{P}^n, \mathbf{Q}^n) = \lim_{n \rightarrow \infty} \sup_{\Theta} | \mathbf{P}(\theta | \mathbf{H}_n) - \mathbf{Q}(\theta | \mathbf{H}_n) | = 0.$$

In other words, the two systems of posterior probabilities *for the parameter*, based on shared evidence, merge together.

Question: What about posterior probability distributions over the algebra generated by the observable events, \mathcal{A} ?

- Recall that the *i.i.d.* Bernoulli statistical model for the data is shared between these two investigators: $(\forall E \in \mathcal{A}) \mathbf{P}(E | \theta) = \mathbf{Q}(E | \theta)$.
- Also, with conjugate priors from the Beta family, the prior probability is positive for each “historical” event H_n . That is, $(\forall H_n, 0 < \theta < 1) \mathbf{P}(H_n | \theta) > 0$.

Moreover, $\mathbf{P}(H_n) = \int_{\Theta} \mathbf{P}(H_n | \theta) d\mathbf{P}(\theta)$. Therefore, $\mathbf{P}(H_n) > 0$, and likewise $\mathbf{Q}(H_n) > 0$.

$$\begin{aligned} \text{Answer: } \mathbf{P}(E | H_n) &= \int_{\Theta} \mathbf{P}(E | \theta, H_n) d\mathbf{P}^n(\theta). \\ &= \int_{\Theta} [\mathbf{P}(E, H_n | \theta) / \mathbf{P}(H_n | \theta)] d\mathbf{P}^n(\theta) \end{aligned}$$

and as \mathbf{P}^n merges with \mathbf{Q}^n for large n ,

$$\approx \int_{\Theta} [\mathbf{P}(E, H_n | \theta) / \mathbf{P}(H_n | \theta)] d\mathbf{Q}^n(\theta)$$

and as the two investigators agree on the statistical model

$$\begin{aligned} &= \int_{\Theta} [\mathbf{Q}(E, H_n | \theta) / \mathbf{Q}(H_n | \theta)] d\mathbf{Q}^n(\theta) \\ &= \int_{\Theta} \mathbf{Q}(E | \theta, H_n) \mathbf{Q}(\theta | H_n) d\mathbf{Q}^n(\theta) \\ &= \mathbf{Q}(E | H_n). \end{aligned}$$

Thus, the two posterior *predictive* distributions (over \mathcal{A}) also merge.

For example, the probability that the next flip lands heads given H_n is:

$$\mathbf{P}(X_{n+1} | H_n) = \mathbf{E}_{\mathbf{P}_n}[\theta] = (\alpha+k) / (\alpha+\beta+n),$$

which for large n ,

$$\approx k/n$$

and by parallel reasoning

$$\approx \mathbf{Q}(X_{n+1} | H_n).$$

Note, that the agreement between $\mathbf{P}(E | H_n)$ and $\mathbf{Q}(E | H_n)$ takes a stronger form for cases when the historical observation H_n precludes E , when $(E \cap H_n) = \emptyset$.

Then,

$$\mathbf{P}(E | H_{n'}) = \mathbf{Q}(E | H_{n'}) = 0 \text{ for all } n' \geq n.$$

Question: What parts of these asymptotic results for the algebra of events \mathcal{A} depends upon the (shared) statistical model?

Answers:

(1) *Asymptotic Certainty* is automatic with the Bayesian framework!

($\forall E \in \mathcal{A}$) with \mathbf{P} -probability 1,

$$\lim_{n \rightarrow \infty} \mathbf{P}^n(E) = \chi(E), \text{ i.e.} \quad (\text{Halmos, 1948})$$

(2) *Asymptotic Consensus* requires only agreement on “null” events.

Assume that ($\forall E \in \mathcal{A}$) $\mathbf{P}(E) = 0$ if and only if $\mathbf{Q}(E) = 0$.

With \mathbf{P} -(or \mathbf{Q} -) probability 1, with respect to \mathcal{A}

$$\lim_{n \rightarrow \infty} \rho(\mathbf{P}^n, \mathbf{Q}^n) = 0. \quad (\text{Blackwell \& Dubins, 1962})$$

However, the statistical model **is** needed for each of the following:

- (1) data reduction
- (2) rates of convergence to certainty
- (3) rates of merging for Bayesian investigators with shared evidence

Next, we explore/review several themes for Bayesian asymptotics:

- A role for *statistical models* in these asymptotic results.
 - symmetry/independence assumptions in these results.
 - data reduction
 - asymptotic Normal inference for these results.

A puzzlement?

We have two investigators (T and J) for our coin-tossing problem. They share the same statistical (*i.i.d.* Bernoulli) model for coin flips, and they have the *same* (conjugate) **Beta** prior for θ .

They collect (shared) evidence by flipping the coin until one says, “Stop.” In fact, they observe the sequence

$(H, H, T, H, T, T, H, H, T, H)$

at which point they both (simultaneously) say “Stop!”

However:

T 's plan was to flip the coin exactly 10 times and stop
and J 's plan was to flip until there were 6 “Heads” and stop.

Exercise: Give the Bayes analysis for T and for J of these data.

Roles for Statistical Models

- Data Reduction and factorization of the likelihood function.
 - *Sufficient* Statistics
 - *Ancillary* Statistics
- Symmetry and Independence assumptions
 - deFinetti's theorem on *exchangeable sequences*
- Properties of *Maximum Likelihood*

Data Reduction Concepts for Statistical Models

Defn: The (dimensional) random variable $Y = \mathbf{g}(X)$ is *sufficient* for the parameter θ (with respect to X) *iff*

$$\mathbf{P}(X | Y, \theta) = \mathbf{P}(X | Y), \text{ independent of } \theta.$$

Theorem: The likelihood for θ given a sufficient (set of) statistic(s) Y is the same as the likelihood for θ given the (dimensional) variable X for which Y is sufficient.

Proof: $\mathbf{P}(x | \theta) = \mathbf{P}(x, y | \theta)$ as $Y = \mathbf{g}(X)$
 $= \mathbf{P}(x | y, \theta) \mathbf{P}(y | \theta)$ multiplication axiom
 $= \mathbf{P}(x | y) \mathbf{P}(y | \theta)$ by sufficiency of Y
 $\propto \mathbf{P}(y | \theta)$

Corollary (Factorization of the likelihood function):

$Y = \mathbf{g}(X)$ is *sufficient* for the parameter θ (with respect to X) *iff*

The likelihood (probability or density) function can be written as the product of two functions of this form:

$$\mathbf{P}(X | \theta) = \mathbf{h}(X) \mathbf{j}(Y, \theta).$$

Recall: $Y = \mathbf{g}(X)$

Example 1 (coin-tossing, again):

$\mathbf{X} = \langle X_1, \dots, X_n \rangle$ are *iid* Bernoulli trials given θ , with $\mathbf{P}(X_1=1|\theta) = \theta$, $0 < \theta < 1$.

Claim: $\mathbf{g}(\mathbf{X}) = \mathbf{Y} = \langle \sum_i X_i, n - \sum_i X_i \rangle$ is a sufficient reduction to the two statistics, #1's = $\sum_i X_i = k$ and #0's = $n - k$ in the sequence \mathbf{X} .

Proof: $\mathbf{P}(\mathbf{x}, \mathbf{y} | \theta) = \mathbf{P}(\mathbf{x} | \theta) = \theta^k (1-\theta)^{n-k}$

$$\mathbf{P}(\mathbf{y} | \theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

Thus $\mathbf{P}(\mathbf{x} | \mathbf{y}, \theta) = \mathbf{P}(\mathbf{x}, \mathbf{y} | \theta) / \mathbf{P}(\mathbf{y} | \theta) = \mathbf{P}(\mathbf{x} | \mathbf{y}) = k!(n-k)!/n!$

That is, $\mathbf{P}(\mathbf{X} | \mathbf{y}, \theta)$ is a discrete, uniform distribution over all sequences H_n that begin with k 1's and $(n-k)$ 0's, independent of θ .

Or, use factorization and note that, alternatively $\langle \bar{X}, n \rangle$ are sufficient for θ as

$$\mathbf{P}(\mathbf{X} | \theta) = \theta^{n\bar{X}} (1-\theta)^{n(1-\bar{X})} = \mathbf{h}(\mathbf{X}) \mathbf{j}(Y, \theta)$$

where $\mathbf{h}(\mathbf{X}) = 1$ and $Y = \bar{X} = \sum_i X_i / n$.

Example 2 (Normal distribution, known variance):

$\mathbf{X} = \langle X_1, \dots, X_n \rangle$ are *iid* normal $N(\mu, 1)$ trials.

Claim: The pair $\langle \bar{X}, n \rangle$ is sufficient for μ .

Proof: Write $\mathbf{p}(X | \mu) =$

$$(2\pi)^{-n/2} \exp(-\sum_i (X_i - \bar{X})^2 / 2) \exp(-n(\mu - \bar{X})^2 / 2),$$

where

$$\underline{(2\pi)^{-n/2} \exp(-\sum_i (X_i - \bar{X})^2 / 2)} \exp(-n(\mu - \bar{X})^2 / 2)$$

↑↑

$h(X)$

↑↑

$j(Y, \theta)$

Defn: The (dimensional) random variable $Y = \mathbf{g}(X)$ is *ancillary* for the parameter θ (with respect to X) *iff*

$$\mathbf{P}(Y | \theta) = \mathbf{P}(Y), \text{ independent of } \theta.$$

Theorem: The likelihood for θ based on an ancillary (set of) statistic(s) Y is *constant*.

Corollary: The likelihood for θ based on X equals the conditional likelihood for θ based on X , *given* Y .

$$\mathbf{P}(x | \theta) = \mathbf{P}(x | y, \theta)$$

Proof: $\mathbf{P}(x | \theta) = \mathbf{P}(x, y | \theta) = \mathbf{P}(x | y, \theta) \mathbf{P}(y | \theta)$
 $\propto \mathbf{P}(x | y, \theta).$

Example 3 (coin-tossing, again):

$\mathbf{X} = \langle X_1, \dots, X_i, \dots \rangle$ are *iid* Bernoulli trials given θ , with $\mathbf{P}(X_1=1|\theta) = \theta$, $0 < \theta < 1$.

$\mathbf{g}(\mathbf{X}) = \mathbf{Y} = \langle \sum_i X_i, n - \sum_i X_i \rangle$ is a sufficient reduction for inference about θ .

Version 3a: The *stopping rule* is sample to a fixed sample size n . Then N (sample size) is ancillary ($\mathbf{P}(N=n) = 1$) and, **given** $N = n$, $\sum_i X_i$ is sufficient!

Moreover, $\mathbf{P}(\sum_i X_i | n, \theta)$ is given by the *Binomial*(n, θ) distribution.

Version 3b: The *stopping rule* is sample to a fixed number of “heads,” say $\sum_i X_i = k$

Then $\sum_i X_i$ (number of heads) is ancillary ($\mathbf{P}(\sum_i X_i = k) = 1$) and, **given** $\sum_i X_i = k$, the number of flips N is sufficient!

Moreover, $\mathbf{P}(N | k, \theta)$ is given by the *Neg-Binomial*(k, θ) distribution.

However, regardless of the stopping rule, in either version, the *pair* $\langle \sum_i X_i, N \rangle$ is sufficient!

Recapitulation of data-reduction principles for statistical models

Sufficiency principle: A sufficient statistic preserves all the relevant information about the parameter that is in the full data set

Ancillarity principle: All the relevant information in the data set about the parameter is contained in the conditional model, given the ancillary statistic.

Likelihood principle: All the relevant information in the data set about the parameter is contained in the likelihood function given the data.

Birnbaum's Theorem: The *Likelihood* principle is equivalent to the conjunction of the *Sufficiency* and *Ancillarity* principles.

Identifying statistical models by symmetry & independence involving observables
(*deFinetti's Theorem*)

Heuristic Example (coin-tossing yet again!): Let $\mathbf{X} = \langle X_1, \dots, X_i, \dots \rangle$ be an infinite sequence of binary trials, with the σ -algebra (\mathbf{A}) of events generated by the observable “historical” events $H_n: \langle x_1, \dots, x_n, \{0,1\}, \{0,1\}, \dots \rangle$.

Defn: Say that a probability \mathbf{P} over \mathbf{A} is:

- *1-exchangeable* if for $\forall(i,j) \mathbf{P}(X_i = 1) = \mathbf{P}(X_j = 1)$

- *2-exchangeable* if $\forall(i_1, i_2, \text{distinct and } j_1, j_2 \text{ distinct})$

$$\mathbf{P}(X_{i_1} = x_1, X_{i_2} = x_2) = \mathbf{P}(X_{j_1} = x_1, X_{j_2} = x_2)$$

- *n-exchangeable* if $\forall(i_1, i_2, \dots, i_n \text{ distinct})$

$\mathbf{P}(X_{i_1} = x_1, X_{i_2} = x_2, \dots, X_{i_n} = x_n)$ does not depend on the n distinct $\langle i_1, i_2, \dots, i_n \rangle$

- *exchangeable* if \mathbf{P} is n -exchangeable for each n ($n = 1, 2, \dots$).

Theorem (deFinetti): \mathbf{P} is exchangeable if and only if \mathbf{P} can be written as

$$\mathbf{P}(E) = \int_{\Theta} \mathbf{P}(E \mid \theta) d\mathbf{Q}(\theta)$$

where

- $\mathbf{P}(A \mid \theta)$ is given by *iid* Bernoulli(θ) trials
- $\mathbf{Q}(\theta)$ is a prior probability distribution over Θ determined uniquely by \mathbf{P} over A .

Thus, one can use the computational benefits of sampling from an *iid* statistical model, “as if” it were true, given suitable exchangeability (symmetry) assumptions involving only the algebra of the observable random variables.

Remarks:

- This important theorem generalizes to cover both discrete and continuous random variables.
- Also, there is version dealing with finite sequences (N-exchangeability).
- For a thorough discussion of all this, see chapter 1 of Mark Schervish’s book, *Theory of Statistics*, 1995. Springer-Verlag.

Data reduction, Fisher-Information, and Maximum Likelihood

Defn.: Score function:
$$\mathbf{S}_X(\theta) = \frac{\partial (\ln \mathbf{p}(X | \theta))}{\partial \theta}$$

Fisher Information (under general conditions)

$$I_X(\theta) = \text{Var}(\mathbf{S}_X(\theta)) = \mathbb{E}\left[-\frac{\partial^2 (\ln \mathbf{p}(X | \theta))}{\partial \theta^2}\right].$$

- Fisher Information is additive for independent data.
- $I_X(\theta) = I_Y(\theta)$ whenever Y is sufficient for θ (with respect to X).
- Fisher Information is a differential form of Kullback-Leiber information.

K-L information inequality

With f and g density functions:

$$\begin{aligned}\int \mathbf{log} [f(\mathbf{x})/g(\mathbf{x})] f(\mathbf{x}) \, d\mathbf{x} &= -\int \mathbf{log} [g(\mathbf{x})/f(\mathbf{x})] f(\mathbf{x}) \, d\mathbf{x} \\ &\geq -\mathbf{log} \left(\int [g(\mathbf{x})/f(\mathbf{x})] f(\mathbf{x}) \, d\mathbf{x} \right) \\ &= -\mathbf{log} \left(\int g(\mathbf{x}) \, d\mathbf{x} \right) \\ &= -\mathbf{log} (1) \\ &= 0\end{aligned}$$

where the inequality follows is by an application of Jensen's inequality for the concave \mathbf{log} function.

The inequality is strict unless almost surely w.r.t. the \mathbf{F} -distribution, $[g(\mathbf{x})/f(\mathbf{x})] = 1$.

Corollary: Unless a statistic t is sufficient for data X , there will be an information (and Fisher Information) gain from adding the full data X to t .

That is, all and only sufficient statistics preserve all the information in the data.

Defn.: Let θ^* denote the argmax of the likelihood function $\mathbf{p}(X | \theta)$,
the *maximum likelihood estimate (MLE)* of the parameter.

Main MLE Theorem (under general regularity conditions on the statistical model):

$$\mathbf{P}(\theta^* | \theta_0) \approx N(\theta_0, [I_X(\theta^*)]^{-1}) = N(\theta_0, [nI_{X_i}(\theta^*)]^{-1})$$

So (under “regularity” conditions) the *MLE*:

- Has an asymptotic Normal distribution.
- Is asymptotically consistent (converges to θ_0).
- Is asymptotically sufficient.

Putting these pieces together we have, under the same “regularity” conditions,
convergence of the posterior to the *mle*:

$\mathbf{P}\left[\frac{(\theta^* - \theta_0)}{[\mathbf{I}_{X_n}(\theta^*)]^{1/2}} \mid X_n\right]$ converges to the standard normal $\mathbf{N}(\mathbf{0},1)$ distribution.