

NONPARAMETRIC BINARY REGRESSION WITH RANDOM COVARIATES

BY

P. DIACONIS* (CAMBRIDGE, MASSACHUSETTS)

AND

D. FREEDMAN** (BERKELEY, CALIFORNIA)

Abstract. The performance of Bayes' estimates is studied under an assumption of conditional exchangeability. More exactly, for each subject in a data set, let ξ be a vector of binary covariates and let η be a binary response variable with $P\{\eta = 1 \mid \xi\} = f(\xi)$. Here, f is an unknown function, to be estimated from the data; the subjects are independent, and the ξ 's are iid uniform. Define a prior distribution on f as $\sum_k w_k \pi_k / \sum_k w_k$, where π_k is uniform on the set of f which only depend on the first k covariates and $w_k > 0$ for infinitely many k . Bayes' estimates are consistent at all f if w_k decreases rapidly as k increase. Otherwise, the estimates are inconsistent at $f \equiv 1/2$.

1. Introduction. This paper studies non-parametric binary regression in a Bayesian context. Let η be a binary response variable and $\xi \in [0, 1]$ a covariate with $P\{\eta = 1 \mid \xi\} = f(\xi)$. The function f is an infinite-dimensional parameter to be estimated from data by Bayesian techniques. At stage n , the data consist of n iid pairs $(\eta_1, \xi(1)), \dots, (\eta_n, \xi(n))$ with $\xi(i)$ uniform over $[0, 1]$. Thus, we are assuming that

(1.1) Given the covariates, the response variables are independent across subjects and $P\{\eta_i = 1 \mid \xi(i)\} = f(\xi(i))$.

The function f is assumed to be measurable.

The main issue to be studied is consistency: does the posterior concentrate near the true f as more data come in? For the class of priors we consider, the answer is generally yes, but not always. For some priors, the Bayes estimates are inconsistent when $f \equiv 1/2$.

We next describe the priors to be considered; these were motivated by an example of de Finetti [3]. Regard a point x in the unit interval as an infinite

* Research partially supported by NSF Grant DMS 86-00235.

** Research partially supported by NSF Grant DMS 92-08677.

sequence of binary digits or "bits," to be denoted by x_1, \dots, x_n . Our priors are "hierarchical" or "nested." We begin with a prior π_k supported on the class of functions f that depend only on the first k covariates, so $\xi_{k+1}, \xi_{k+2}, \dots$ do not matter in (1.1).

From the perspective of π_k , $P\{\eta_i = 1 \mid \xi(i)\}$ depends only on $\xi_1(i), \dots, \xi_k(i)$. Let C_k be the set of strings of 0's and 1's of length k . So, the prior π_k is defined by the joint distribution it assigns to the 2^k success probabilities $\theta_s, s \in C_k$. Here, $\theta_s = P\{\eta = 1 \mid \xi\}$ for subjects with $(\xi_1, \dots, \xi_k) = s$. One simple choice for π_k takes θ_s independent and uniform over $[0, 1]$, as s varies over C_k ; that is the example to keep in mind.

We treat k as an unknown "hyperparameter," putting a prior weight w_k on k . Thus, our prior is of the form

$$(1.2a) \quad \pi = \sum_{k=0}^{\infty} w_k \pi_k / \sum_{k=0}^{\infty} w_k,$$

where

$$(1.2b) \quad w_k > 0 \text{ for infinitely many } k \text{ and } \sum_{k=0}^{\infty} w_k < \infty.$$

Our results apply to a more general class of " Γ -uniform" π_k , to be defined now. Fix $0 < b \leq B < \infty$ and a finite subset F of $(0, 1)$. Consider the class Γ of all densities γ on $[0, 1]$ with $b \leq \gamma \leq B$. Let $g_s = \int \theta \gamma_s(\theta) d\theta$. By assumption, the g_s all lie in the finite subset F of $(0, 1)$, which was given a priori. We require π_k to make the 2^k success probabilities $\theta_s, s \in C_k$, independent, with densities in the class Γ . Furthermore, we require the choices to fit together as k varies, in the following manner: There is a continuous function $g_{\infty}(x)$, which takes values in F , such that for all $n > n_0$ and all $s \in C_n$, $g_s = g_{\infty}(x)$ provided the first n bits of x agree with s . (Of course, if a continuous function on sequence space takes only finitely many values, it must be piecewise constant.) This completes the definition of Γ -uniformity. If $b = B = 1$ and $F = \{1/2\}$, we get back to the uniform priors.

To define consistency, we need to topologize the parameter space of functions f . Let $C_{\infty} = \{0, 1\}^{\infty}$ be the space of sequences of 0's and 1's, so f maps C_{∞} into $[0, 1]$. Write λ^{∞} for coin-tossing measure on C_{∞} , which makes the bits of $x \in C_{\infty}$ independent, each being 1 with probability $1/2$ and 0 with probability $1/2$. By definition, the parameter space Θ consists of all measurable functions from C_{∞} to $[0, 1]$; functions which are equal a.e. are identified. Put the L_2 -metric on Θ . (The same topology is given by L_p for $p \geq 1$, or convergence in measure.) A typical neighborhood $N(f, \delta, \varepsilon)$ of f is defined as follows:

$$(1.3) \quad \text{If } f \in \Theta \text{ and } \delta, \varepsilon > 0, \text{ let } N(f, \delta, \varepsilon) \text{ be the set of } h \in \Theta \text{ with}$$

$$\lambda^{\infty} \{x: x \in C_{\infty} \text{ and } |h(x) - f(x)| \leq \varepsilon\} \geq 1 - \delta.$$

If π is a prior probability on Θ , the posterior probability $\tilde{\pi}_n$ on Θ is the conditional law of f , given the data at stage n ; this will be computed explicitly in Section 2. The prior π is *consistent at f* if $\tilde{\pi}_n\{N(f, \delta, \varepsilon)\} \rightarrow 1$ almost surely as $n \rightarrow \infty$. Here, the data are generated in accordance with (1.1), so f is the "true" regression function.

At stage n , there are n subjects, indexed by $i = 1, \dots, n$. Each subject i has a response variable $\eta_i = \eta(i)$ and an infinite sequence of covariate bits

$$\xi(i) = \xi_1(i), \xi_2(i), \dots$$

In addition to (1.1), we are assuming

(1.4) The $\xi(i)$ are independent and have a common uniform distribution.

The main results can now be stated.

(1.5) THEOREM. *Suppose (1.1) and (1.4). Moreover, suppose that the π_k are Γ -uniform, the prior π is hierarchical in the sense of (1.2), and $f \neq g_\infty$. Then π is consistent.*

(1.6) THEOREM. *Suppose (1.1) and (1.4). Moreover, suppose that the π_k are Γ -uniform, the prior π is hierarchical in the sense of (1.2), and $f \equiv g_\infty$. Let l be the smallest k with $w_k > 0$. Write $\exp(x) = e^x$.*

(a) *Suppose $\sum_{k=n}^{\infty} w_k < \exp(-\frac{1}{4}(\log 2)n2^l - \delta_0 n2^l)$ for all large n , for some $\delta_0 > 0$. Then π is consistent at f .*

(b) *Suppose $\sum_{k=n}^{\infty} w_k > \exp(-\frac{1}{4}(\log 2)n2^l + \delta_0 n2^l)$ for infinitely many n , for some $\delta_0 > 0$. Then π is inconsistent at f .*

Theorems (1.5) and (1.6) show that our Bayes estimates are consistent, provided the weights w_k fall off rapidly. For example, suppose $b = B = 1$, so that π_k makes the success probabilities iid uniform. If $w_k = 1/2^k$, then the Bayes estimates are consistent at all f . On the other hand, if $w_k = 1/(k+1)^2$, the estimates are inconsistent at $f \equiv 1/2$.

The present paper extends results in Diaconis and Freedman [8]. That paper studies the model (1.1), but with a different sampling design from (1.4). The $\xi(i)$ were taken to be balanced, so that at stage n all $s \in C_n$ occur once and only once. The balance condition eliminates some annoying inhomogeneities that have to be dealt with here by Poissonization. However, the critical rate for consistency depends on the sampling plan: compare (1.6) with Theorems (8) and (9) in [8].

In [8] we review the relationship between consistency of Bayes estimates and rules for model selection as well as sieves and orthogonal series estimation. Also see Diaconis and Freedman [9], where we try to identify the root cause of the inconsistency in (1.6b) and suggest that consistency will hold fairly generally; we also suggest that Bayesian methods will generally lead to correct estimates of the order when the model is finite.

The balance of this paper is organized as follows. Section 2 computes the posterior. Section 3 has some preliminary estimates, including large-deviation results for balls dropped at random into boxes. A proof of Theorem (1.5) is given in Section 4, and Theorem (1.6) will be proved in Section 5. The arguments are modifications of those in [8].

Consistency of Bayes estimates is a classical problem going back to Laplace [16] and the present paper uses methods introduced by him. A recent survey, with emphasis on infinite-dimensional problems, will be found in [6]. The combinatorial literature on dropping balls at random into boxes is surveyed by Kolchin et al. [15]. A recent treatment using Stein's method is in [1]. Entry-points to the large-deviation literature are [11] or [17]; papers [4] or [5] are more recent. The related topic of boundary crossings is surveyed in [18]. Poissonization can be traced back to [12]. A modern reference is [14].

2. Computing the posterior. Let Ω be an underlying probability space, on which the response variables $\eta(i)$ and covariates $\xi_j(i)$ are defined. Recall that $f \in \Theta$ maps C_∞ to $[0, 1]$. For $f \in \Theta$, let P_f be the probability on Ω which makes the response variables and covariates distributed so that (1.1) and (1.4) hold. The dependence between the data at stage n and stage $n+1$ is simple: there is one extra subject with covariate sequence $\xi(n+1)$. The joint distribution across n 's will matter for some of the arguments here, as opposed to [8].

Let

$$(2.1) \quad f_k(x) = E\{f | x_1, \dots, x_k\} = \int_{C_\infty} f(x_1, \dots, x_k, y) \lambda^\infty(dy)$$

and write $f_k(s)$ for $f_k(x)$ when $s \in C_k$ and $x_1 = s_1, \dots, x_k = s_k$.

For now, fix n and k . For $s \in C_k$, let N_s be the number of subjects $i = 1, \dots, n$ such that $\xi_j(i) = s_j$ for $j = 1, \dots, k$. In other words, N_s is the number of subjects $i = 1, \dots, n$ whose first k covariates are given by s . Of course, N_s is random; that is the new technical difficulty. Let X_s be the number of successes among subjects whose covariate sequence begins with s . More formally, $\eta(i)$ is the response for subject i , and

$$(2.2) \quad X_s = \sum_{i=1}^n \{\eta(i) : \xi_j(i) = s_j \text{ for } i = 1, \dots, n\}.$$

Write $\text{bin}(m, p)$ for the binomial distribution with m trials and success probability p .

(2.3) LEMMA. Assume (1.4). With respect to P_f :

(a) $\{N_s : s \in C_k\}$ is distributed like the result of dropping n balls at random into 2^k cells.

(b) Given $\{N_s : s \in C_k\}$, the random variables X_s are independent as s ranges over C_k , each being $\text{bin}[N_s, f_k(s)]$.

As usual, π_k can be extended to a probability on $\Theta \times \Omega$ by the formula

$$\pi_k(A \times B) = \int_A P_f \{B\} \pi_k \{df\},$$

where A is a measurable subset of Θ , and B is a measurable subset of Ω . The proofs of Lemmas (2.3) and (2.4) are omitted as routine. In (2.4) and similar contexts, π_k is viewed as a probability on $\Theta \times \Omega$.

(2.4) LEMMA. Suppose π_k is Γ -uniform. With respect to π_k , the N_s have the ball-dropping distribution given by (2.3). Given $\{N_s: s \in C_k\}$, the pairs (θ_s, X_s) are independent as s ranges over C_k . The parameter θ_s has density $\gamma_s \in \Gamma$. Given N_s , and θ_s , the number of successes X_s is $\text{bin}(N_s, \theta_s)$.

For $\gamma \in \Gamma$, $m = 0, 1, 2, \dots$, and $j = 0, 1, \dots, m$, let

$$(2.5a) \quad \gamma(m, j, \cdot): \theta \rightarrow \frac{\theta^j (1-\theta)^{m-j} \gamma(\theta)}{\phi(m, j, \gamma)},$$

where the normalizing constant is

$$(2.5b) \quad \phi(m, j, \gamma) = \int_0^1 \theta^j (1-\theta)^{m-j} \gamma(\theta) d\theta.$$

In particular, $\phi(0, 0, \gamma) = 1$ and $\gamma(0, 0, \cdot) = \gamma(\cdot)$.

Let $\tilde{\pi}_{k,n}$ be the posterior distribution of f , computed relative to π_k , given the data from a design of order n .

(2.6) LEMMA. Suppose π_k is Γ -uniform. According to the posterior $\tilde{\pi}_{k,n}$, the success probabilities θ_s are independent as s ranges over C_k , and θ_s has density $\gamma_s(N_s, X_s, \cdot)$ with respect to Lebesgue measure on $[0, 1]$.

To compute the posterior relative to π , the π_k -predictive probability of the data is needed. To set up the notation, recall the normalizing constant ϕ from (2.5b). Let

$$(2.7) \quad \varrho_{k,n} = \prod_{s \in C_k} \phi(N_s, X_s, \gamma_s).$$

If $N_s = 0$, the corresponding factor in $\varrho_{k,n}$ is taken as 1. By (2.4), $\varrho_{k,n}$ is the π_k -predictive probability of the data given $\{N_s\}$.

Turn now to the posterior $\tilde{\pi}_n$, computed relative to π . Informally, the "theory index" k in (1.2) is a parameter which has a posterior distribution relative to π . Let

$$(2.8) \quad \tilde{w}_{k,n} = w_k \varrho_{k,n}.$$

Now, $\pi_k \{\text{data}\} / \pi \{\text{data}\} = \tilde{w}_{k,n} / \sum_{k=0}^{\infty} \tilde{w}_{k,n}$. So

$$(2.9) \quad \tilde{\pi}_n(k) = \tilde{w}_{k,n} / \sum_{k=0}^{\infty} \tilde{w}_{k,n}.$$

(2.10) LEMMA. Suppose π is hierarchical in the sense of (1.2), and the π_k are Γ -uniform. Given the data from a design of order n , the posterior is

$$\tilde{\pi}_n = \sum_{k=0}^{\infty} \tilde{w}_{k,n} \tilde{\pi}_{k,n} / \sum_{k=0}^{\infty} \tilde{w}_{k,n}.$$

The proof is omitted as routine. Of course, $\tilde{\pi}_n$ can be written as

$$\sum_{k=0}^{\infty} \tilde{\pi}_n(k) \tilde{\pi}_{k,n}.$$

3. Some estimates.

(3.1) LEMMA. Let $0 \leq p \leq 1$. Let X be $\text{bin}(m, p)$ and $Y = (X - mp)^2/m$. If $m = 0$, or $m > 0$ but $p = 0$ or 1 , let $Y = 0$. Then:

- (a) $P\{X \leq mp - \sqrt{mx}\} < \exp(-\frac{1}{2}x)$ for all $x > 0$.
- (b) $P\{X \geq mp + \sqrt{mx}\} < \exp(-\frac{1}{2}x)$ for all $x > 0$.
- (c) Y is stochastically smaller than $\chi^2_2 + 2\log 2$.

Proof. Suppose $m > 0$ and $0 < p < 1$. Claim (a) follows from Bernstein's inequality. For example, use (4) in [13] to see that

$$P\{X \leq mp - \sqrt{mx}\} < \exp\left(-\frac{1}{2p}x\right) < \exp\left(-\frac{1}{2}x\right).$$

To get claim (b), write $q = 1 - p$, and observe that $X \geq mp + \sqrt{mx}$ iff $(m - X) \leq mq - \sqrt{mx}$. Now use (a). For (c),

$$P\{Y \geq x\} < 2\exp(-\frac{1}{2}x). \quad \blacksquare$$

(3.2) LEMMA. Suppose the random variable ξ has a Laplace transform for $h < h_0$, where h_0 is positive. Let \mathcal{K} be the class of random variables Y for which $E\{[Y - EY]^j\} \leq E\{\xi^j\}$ for $j \geq 2$. Then there are positive finite σ^2 and h_1 , depending only on ξ , such that

$$P\left\{\sum_{i=1}^m Y_i \geq \sum_{i=1}^m E\{Y_i\} + y\right\} < \exp(-y^2/2\sigma^2 m)$$

provided the Y_i 's are independent, $Y_i \in \mathcal{K}$ for all i , and $0 < y \leq h_1 m$.

Note. This lemma is set up to give one-sided bounds. In some cases, of course, it can also be applied to $\{-Y\}$. Then a lower bound can be obtained in the same way with a slightly smaller σ . More detailed results can be obtained by matching variances or Esscher tilting, but these refinements will not be needed here. See [11], Section XVI.6.

Proof. Assume without real loss of generality that $EY = 0$. Let

$$\phi_Y(h) = E\{e^{hY}\} = 1 + \sum_{j=2}^{\infty} E\{Y^j\} h^j/j!.$$

The sum is bounded above by

$$\sum_{j=2}^{\infty} h^j E\{\xi^j\}/j! < \frac{1}{2}\sigma^2 h^2 \quad \text{for } 0 < h \leq h',$$

where σ^2 is a suitable positive finite number, slightly larger than the second moment of ξ . For $0 < h \leq h'$,

$$\phi_Y(h) < 1 + \frac{1}{2}\sigma^2 h^2 \quad \text{and} \quad \log \phi_Y(h) < \frac{1}{2}\sigma^2 h^2.$$

The constants σ^2 and h' depend on ξ , not Y or h .

We are assuming that $E\{Y_i\} = 0$. Chebyshev's inequality can be applied to bound $P\{\exp(h(Y_1 + \dots + Y_m)) \geq e^{hy}\}$:

$$\log P\{Y_1 + \dots + Y_m \geq y\} \leq -hy + m \log \phi_Y(h) < -hy + \frac{1}{2}m\sigma^2 h^2.$$

Put $h = y/\sigma^2 m$. We require $h \leq h'$, i.e., $y \leq h'\sigma^2 m$: set $h_1 = h'\sigma^2$. ■

(3.3) LEMMA. Suppose $|U|$ is stochastically smaller than V . Then $|U - EU|$ is stochastically smaller than $V + EV$.

(3.4) COROLLARY. Let n_i be non-negative integers and $0 \leq p_i \leq 1$. Let X_i be independent bin(n_i, p_i) and $Y_i = (X_i - n_i p_i)^2/n_i$. Then there are universal positive constants σ^2 and h_1 such that

$$P\{Y_1 + \dots + Y_m \geq \sum_{i=1}^m p_i(1-p_i) + y\} < \exp(-y^2/2\sigma^2 m)$$

provided $0 < y \leq h_1 m$.

Proof. Combine (3.1)–(3.3). ■

(3.5) LEMMA. Let N_λ be Pois(λ), i.e., Poisson with parameter λ . If $N_\lambda = 0$, let $\log(N_\lambda) = 0$. Let $z > 0$. Then:

(a) $P\{\sqrt{\lambda}(\log N_\lambda - \log \lambda) \geq z\} < \exp(-\frac{1}{2}z^2)$ for all $z > 0$.

(b) $P\{\sqrt{\lambda}(\log N_\lambda - \log \lambda) \leq -z\} < \exp(-\frac{1}{2}z^2 [(1 - e^{-\varepsilon})/\varepsilon]^2)$

provided $0 < z \leq \varepsilon\sqrt{\lambda}$. If $\varepsilon = 1/2$, then an upper bound is $\exp(-z^2/4)$.

Proof. This follows from Bernstein's inequality: see (4) in [13]. Some auxiliary calculations are needed to estimate the function in (9) of that paper. We claim:

(3.6) $u \rightarrow (e^u - 1)^2/u^2 e^u$ is strictly convex with a minimum at $u = 0$.

Indeed, the function in (3.6) is $((e^{u/2} - e^{-u/2})/u)^2$, which is readily expanded in even powers of u , with positive coefficients.

(3.7) $\lambda(e^{z/\sqrt{\lambda}} - 1)^2/e^{z\sqrt{\lambda}} > z^2$ for $\lambda > 0$ and $z > 0$.

(3.8) $u \rightarrow (1 - e^{-u})/u$ is strictly decreasing for $u > 0$.

(3.9) $\lambda(1 - e^{-z/\sqrt{\lambda}})^2 > z^2((1 - e^{-\varepsilon})/\varepsilon)^2$ for $0 < z < \varepsilon\sqrt{\lambda}$.

(3.10) COROLLARY. Let $N'_\lambda = N_\lambda$ if $N_\lambda > \lambda e^{-1/2}$, else let $N'_\lambda = \lambda e^{-1/2}$. Let

$$Z_\lambda = \sqrt{\lambda} (\log N'_\lambda - \log \lambda).$$

Then Z_λ^2 is stochastically smaller than $4 \log 2 + 2\chi_2^2$.

(3.11) COROLLARY. There are finite positive constants σ^2 and h_1 such that

$$P \left\{ \left| \sum_{i=1}^m (Z_i - E\{Z_i\}) \right| \geq y \right\} < 2 \exp(-y^2/2\sigma^2 m)$$

provided the Z_i are independent, each Z_i is distributed as Z_{λ_i} in (3.10), and $0 < y \leq h_1 m$.

(3.12) LEMMA. $\lim_{\lambda \rightarrow \infty} E\{\log N'_\lambda\}/\log \lambda = \lim_{\lambda \rightarrow \infty} E\{\log N_\lambda\}/\log \lambda = 1$.

Note. As $\lambda \rightarrow \infty$, the law of Z_λ tends to the standard normal. The bound in (3.5b) can be improved, but there is mass $P\{N_\lambda = 0\} = e^{-\lambda}$ at $z = -\sqrt{\lambda} \log \lambda$; no upper bound of the form $\exp(-\delta z^2)$ can be valid for large λ .

Results (3.13)–(3.16) are familiar, but are included for ease of reference. The elementary proof of (3.13) is omitted.

(3.13) LEMMA. Let f be a convex function. Let $a, b > 0$ and let the random variable X_{ab} take values $-a$ or b and $E\{X_{ab}\} = 0$. Then $E\{f(X_{ab})\}$ increases with b for fixed a ; likewise, $E\{f(X_{ab})\}$ increases with a for fixed b .

(3.14) LEMMA. Let f be a convex function. Fix A, B and μ with $-\infty < A < \mu < B < \infty$. Let \mathcal{X} be the class of random variables X such that $A \leq X \leq B$ and $E\{X\} = \mu$. Let $\xi \in \mathcal{X}$ take only the values A, B and $E\{\xi\} = \mu$. Then

$$E\{f(X)\} \leq E\{f(\xi)\}.$$

Proof. Assume without loss of generality that $\mu = 0$. The extreme X have two-point distributions and (3.13) applies. ■

(3.15) COROLLARY. Let f be convex and increasing. Fix L and ε positive and finite. Let \mathcal{X} be the class of random variables X such that $|X| \leq L$ and $E\{X\} \leq -\varepsilon$. Let $\xi \in \mathcal{X}$ take only the values $\pm L$ and $E\{\xi\} = -\varepsilon$. Then $E\{f(X)\} \leq E\{f(\xi)\}$.

(3.16) LEMMA. Define \mathcal{X} as in (3.15). There is a q with $0 < q < 1$, depending only on L and ε , such that: for independent $X_i \in \mathcal{X}$ and $y > 0$,

$$P \left\{ \sum_{i=1}^m X_i \geq y \text{ for some } m \right\} < q^y.$$

Proof. Define ξ as in (3.15). De Moivre solved the gambler's ruin problem by finding the unique $r > 1$ with $E\{r^\xi\} = 1$. Continuing his argument, let $S(m) = \sum_{i=1}^m X_i$. Then $r^{S(m)}$ is an expectation-decreasing martingale, which can be stopped at the crossing time; take $q = 1/r$. ■

Remark. Lemma (3.16) is easily extended to partial sums of variables X_i such that the conditional law of X_i gives the past falls in \mathcal{K} ; see [10], p. 164.

Lemmas (3.17) and (3.18) are elementary, and proofs are omitted.

(3.17) LEMMA. Let j be a non-negative integer, and x be a positive real number. Let $f_j(x) = \sum_{i=j}^{\infty} x^i/i!$. Then $f_j(x)/x^j$ is continuous, convex, and strictly increasing on $(0, \infty)$, with a limit of $1/j!$ as x decreases to 0.

(3.18) LEMMA. Let m be a positive integer and $0 < p < 1$. Let X be $\text{bin}(m, p)$, and let j be a non-negative integer. Then

- (a) $P\{X = j\} < (j!)^{-1} (mp)^j$;
- (b) $P\{X \geq j\} < f_j(mp)$.

(3.19) LEMMA. Assume (1.4). Fix $c > 5/3$. Then almost surely, for all sufficiently large n , for all $k > c \log_2 n$, there are no $s \in C_k$ with $N_s \geq 4$.

Proof. By (2.3a), N_s is $\text{bin}(n, 1/2^k)$. Write $\lambda = n/2^k = E\{N_s\}$. So $\lambda < 1/n^{c-1}$. By (3.17) and (3.18), $P\{N_s \geq 4\} < C\lambda^4$, where C is a suitable positive constant (a bit larger than $1/4!$). The expected number of s with $N_s \geq 4$ is then smaller than $C2^k\lambda^4 = Cn\lambda^3 < C/n^{3c-4}$. The chance of having at least one box with $N_s \geq 4$ is also smaller than C/n^{3c-4} , by Chebyshev's inequality. Since $3c > 5$, the Borel-Cantelli lemma implies that for all sufficiently large n , for k the least integer exceeding $c \log_2 n$, there are no $s \in C_k$ with $N_s \geq 3$. Finally, for n fixed, $|\{s: s \in C_k \text{ and } N_s \geq 4\}|$ is decreasing as k increases. ■

Note. We write $|S|$ for the cardinality of a set S .

(3.20) LEMMA. Assume (1.4). Fix $c > 7/4$. Then almost surely, for all sufficiently large n , for all $k > c \log_2 n$,

- (i) there are no $s \in C_k$ with $N_s \geq 4$, and
- (ii) there is at most one $s \in C_k$ with $N_s = 3$.

Proof. (i) follows from (3.19), since $7/4 > 5/3$. For (ii), let Q_k be the event that $N_s = 3$ for two or more $s \in C_k$. Thus,

$$Q_k = \bigcup \{N_s = 3 \text{ and } N_t = 3 \mid s, t \in C_k \text{ and } s \neq t\}$$

and

$$P\{Q_k\} = \binom{2^k}{2} P_f\{N_s = 3 \text{ and } N_t = 3\}.$$

Now $P_f\{N_s = 3\} < \lambda^3/6$ by (3.18a). Given $\{N_s = 3\}$, N_t is $\text{bin}[n-3, 1/(2^k-1)]$, so

$$P\{N_t = 3 \mid N_s = 3\} < \frac{1}{6} \left(\frac{n-3}{2^k-1} \right)^3 < \lambda^3/6.$$

Thus

$$P\{Q_k\} < \frac{1}{72} 2^{2k} \lambda^6 = \frac{1}{72} n^2 \lambda^4 \leq 1/(72n^{4c-6}).$$

The balance of the argument is omitted, as similar to (3.19). ■

(3.21) LEMMA. Assume (1.4). Fix $c > 3$. Then almost surely, for all sufficiently large n , for all $k > c \log_2 n$, there are no $s \in C_k$ with $N_s \geq 2$.

Proof. Only the minimal k needs to be considered. Now $P_f\{N_s \geq 2\} < C\lambda^2$ by (3.17) and (3.18), so the expected number of $s \in C_k$ with $N_s \geq 2$ is at most $C2^k \lambda^2 = Cn^2/2^k < C/n^{c-2}$. Since $c > 3$, the Borel-Cantelli lemma completes the proof. ■

Lemmas 3.19-3.21 involve the dependence structure of the ball-dropping process as k and n vary. The next result does not. Consider dropping n balls at random into b boxes, where n is much smaller than b : in the case of interest, b is of order $n^2/\log n$. Let $\lambda = n/b$, the expected number of balls in each box.

(3.22) DEFINITION. Let $|M|$ be the number of multiply-occupied boxes, and T the total number of balls in the set of multiply-occupied boxes. Let $S_n = T - |M|$ with $S_0 = 0$. Let p_j be the conditional probability that ball j drops into a previously-occupied box, given the results of dropping the first $j-1$ balls.

Clearly, $S_n \leq n-1$, where the bound is sharp; $n-S_n$ is the number of occupied cells; $S_1 = 0$; and for $n \geq 2$, $S_n = \sum_{j=2}^n X_j$, where X_j is 1 if ball j drops into a previously-occupied box, else X_j is 0. Recall p_j from (3.22). Of course, p_j is itself a random variable, and

$$(3.23) \quad p_j = (j-1-S_{j-1})/b.$$

(3.24) LEMMA. Let $\mu = n(n-1)/2b$.

(a) If $0 < \delta < 1$, then $P\{S_n \geq (1+\delta)\mu\} < \exp(-\delta^2\mu/4)$.

(b) Suppose $0 < \delta < 1$, and $n/b < \delta/2$. Then

$$P\{S_n \leq (1-\delta)\mu\} < \exp(-\delta^2\mu/8).$$

Proof. Claim (a) is Bernstein's inequality: see, e.g., (4) in [13], noting that $\sum_{j=1}^n p_j \leq \mu$.

Claim (b) is similar. Indeed,

$$P\left\{S_n \leq -\frac{\delta}{2}\mu + \sum_{j=1}^n p_j\right\} < \exp(-\delta^2\mu/8).$$

Clearly, $\mu - nS_n/b \leq \sum_{j=1}^n p_j$, because S_j increases with j . So

$$\left\{S_n \leq \left(1 - \frac{\delta}{2}\right)\mu - \frac{nS_n}{b}\right\} \subset \left\{S_n \leq -\frac{\delta}{2}\mu + \sum_{j=1}^n p_j\right\}.$$

Furthermore, $S_n \leq (1 - \delta/2)\mu - nS_n/b$ iff $S_n \leq \alpha\mu$, where

$$\alpha = \frac{1 - \delta/2}{1 + n/b} > 1 - \delta.$$

Therefore, $\{S_n \leq (1 - \delta)\mu\} \subset \{S_n \leq \alpha\mu\}$. ■

Note. The argument shows S_n to be stochastically smaller than $\sum_{j=2}^n Y_j$, where the Y_j are independent 0-1 valued random variables, and $P\{Y_j = 1\} = (j-1)/b$.

(3.25) LEMMA. Fix $j \geq 0$. Let N_λ be Poisson, but conditioned to be j or more. Then N_λ is stochastically increasing with λ .

Proof. Let $f_j(\lambda) = \sum_{k=j}^{\infty} \lambda^k/k!$. If $i > j$, we claim that $f_i(\lambda)/f_j(\lambda)$ increases with λ . This comes down to showing

$$(3.26) \quad f'_i(\lambda)/f_i(\lambda) > f'_j(\lambda)/f_j(\lambda).$$

However,

$$f'_i(\lambda) = \frac{1}{(i-1)!} \lambda^{i-1} + f_i(\lambda).$$

So (3.26) in turn reduces to

$$(3.27) \quad \sum_{k=j}^{\infty} \lambda^{k-j}/k! > \sum_{k=i}^{\infty} \lambda^{k-i}(i-1)!/k!,$$

which holds term by term. ■

4. The proof of Theorem (1.5). This is proved like (8) in [8]; only the main points are given. Zones are defined in terms of positive integers K_i to be chosen later.

Early zone: $0 \leq k \leq K_1$.

Lower midzone: $K_1 \leq k \leq \log_2 \log n + K_2$.

Upper midzone: $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$.

Endzone: $\log_2 n - K_3 \leq k \leq \log_2 n + K_4$.

High zone: $\log_2 n + K_4 \leq k$.

The endzone and high zone have negligible posterior mass; the early zone is negligible too unless $f = f_k$ for some k . Almost surely, for all large n , for all k in the midzone, for most $s \in C_k$, N_s is large and the MLE $\hat{p}_s = X_s/N_s$ is close to $f_k(s)$. Of course, the latter tends to f : see (2.1). Finally, the posterior piles up around the MLE by [7]. We turn to details; Lemmas (4.2)–(4.4) do most of the work for the midzone.

$$(4.1) \quad \text{Let } \lambda = n/2^k, \text{ so } E\{N_s\} = \lambda.$$

(4.2) LEMMA. Assume (1.1) and (1.4). Fix any positive integer K . Then almost surely $[P_f]$, for all sufficiently large n , for all k with $0 \leq k \leq \log_2 \log n + K$ and all $s \in C_k$:

$$(a) N_s > n/2^{k+1} \geq n/(2^{K+1} \log n).$$

$$(b) |\hat{p}_s - f_k(s)| < (2\sqrt{2^{K+1} \log n})/\sqrt{n}.$$

Proof. Claim (a). By (2.3a), N_s is bin($n, 1/2^k$). Abbreviate $C = 1/2^{2K+3}$. By Bernstein's inequality (3.1a),

$$P\{N_s \leq \lambda/2\} < \exp(-\lambda^2/8n) \leq \exp[-Cn/(\log n)^2].$$

The number of strings $s \in C_k$ with $0 \leq k \leq \log_2 \log n + K$ is

$$\sum_{k=0}^{\log_2 \log n + K} 2^k \leq 2^{K+1} \log n$$

and

$$\sum_{n=1}^{\infty} (\log n) \exp[-Cn/(\log n)^2] < \infty.$$

The Borel-Cantelli lemma completes the proof of (a).

The proof of (b) is similar. Indeed, by (3.1),

$$P_f\{|X_s - N_s f_k(s)| \geq \sqrt{N_s} \cdot 2\sqrt{\log n}\} < 2\exp(-2\log n) = 2/n^2. \blacksquare$$

(4.3) LEMMA. Assume (1.1) and (1.4). Fix any large finite M and small positive δ . Then there are positive integers K_2 and K_3 (depending on M and δ) such that: almost surely $[P_f]$, for all sufficiently large n , for all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$, for all but $\delta 2^k$ strings $s \in C_k$, $N_s > M$.

Proof. The argument is by Poissonization. For now, fix k . Let N_s^* be iid Pois(λ) as s varies over C_k . Thus, $\{N_s\}$ is distributed as $\{N_s^*\}$, given that $\{\sum_s N_s^* = n\}$. The conditioning event has probability asymptotic to $1/\sqrt{2\pi n}$. Choose K_3 so large that

$$P\{\text{Pois}(2^{K_3}) \leq M\} < \delta/2.$$

The chance that $\delta 2^k$ more of the $s \in C_k$ have $N_s^* \leq M$ is bounded above by $\exp(-\delta^2 2^k/8)$. This follows from Bernstein's inequality (3.1b); also see (3.25). Now $2^k \geq 2^{K_2} \log n$. Choose K_2 so large that $C = 2^{K_2} \delta^2/8 > 1.5$. There are fewer than $\log_2 n$ theories k to consider, and

$$\sum (\log_2 n) \sqrt{n}/n^C < \infty.$$

The Borel-Cantelli lemma completes the proof. \blacksquare

(4.4) LEMMA. Assume (1.1) and (1.4). Fix δ and ε positive but small. Then there are positive integers K_2 and K_3 such that: almost surely $[P_f]$, for all sufficiently

large n , for all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$, for all but $\delta 2^{k+1}$ strings $s \in C_k$, $|\hat{p}_s - f_k(s)| < \varepsilon$.

Proof. By Chebyshev's inequality, if X is $\text{bin}(m, p)$, then

$$P\{|X - mp| \geq \varepsilon m\} \leq 1/(4\varepsilon^2 m).$$

Choose M finite but so large that $1/(4\varepsilon^2 M) < \delta/2$. By (4.3), apart from $\delta 2^k$ strings $s \in C_k$, $N_s > M$. Given $\{N_s\}$, the X_s are independent $\text{bin}[N_s, f_k(s)]$ random variables. Bernstein's inequality — with no Poissonization needed — completes the argument, as in (4.3): There are other $\delta 2^k$ exceptional strings, and setting them aside, we obtain $|\hat{p}_s - f_k(s)| < \varepsilon$. ■

The early zone: $k \leq K_1$.

Let

$$(4.5) \quad L_{k,n} = n^{-1} \log \varrho_{k,n} = n^{-1} \sum_{s \in C_k} \log \phi(N_s, X_s, \gamma_s).$$

We also need the entropy function:

$$(4.6) \quad H(p) = \begin{cases} p \log p + (1-p) \log(1-p) & \text{for } 0 < p < 1, \\ 0 & \text{for } p = 0 \text{ or } 1. \end{cases}$$

(4.7) LEMMA. Suppose (1.1) and (1.4). Suppose the π_k are Γ -uniform. Fix k . Then

$$\lim_{n \rightarrow \infty} L_{k,n} = \int H(f_k) d\lambda^\infty \text{ almost surely } [P_f].$$

Proof. This is like (4.12) in [8]. Since k is fixed, C_k is finite. We have $N_s \approx n/2^k$ almost surely by the ordinary strong law: see (2.3a). And $\hat{p}_s \rightarrow f_k(s)$ by the strong law or (4.2b). By (3.2)–(3.3c) in [8],

$$n^{-1} \log \phi(N_s, X_s, \gamma_s) \rightarrow 2^{-k} H[f_k(s)] \text{ a.s. } \blacksquare$$

The endzone: $\log_2 n - K_3 \leq k \leq \log_2 n + K_4$.

(4.8) LEMMA. Suppose (1.1) and (1.4). Suppose the π_k are Γ -uniform. Fix any positive integers K_3 and K_4 . Then there is a positive $\varrho < 1$, a finite positive constant A , and a small positive δ (all depending on K_3 and K_4) such that, for all n , for all k with $\log_2 n - K_3 \leq k \leq \log_2 n + K_4$,

$$P_f\{L_{k,n} \geq \int H(f_k) d\lambda^\infty - \delta\} < A \sqrt{n} \varrho^n.$$

Proof. The argument proceeds by Poissonization, as in (4.3). For the moment, fix k . Recall that $\lambda = n/2^k$.

(4.9a) Let N_s^* be iid $\text{Pois}(\lambda)$ for $s \in C_k$.

(4.9b) Given $\{N_s^*\}$, let the X_s^* be independent $\text{bin}[N_s^*, f_k(s)]$.

Let

$$(4.9c) \quad Y_s^* = \log \int_0^1 \theta^{X_s^*} (1-\theta)^{N_s^* - X_s^*} \gamma_s(\theta) d\theta.$$

It suffices to prove

$$(4.10) \quad P \left\{ n^{-1} \sum_{s \in C_k} Y_s^* \geq \int H(f_k) - \delta \right\} < \varrho^n.$$

Choose L^* with $2 \leq L^* < \infty$. We claim:

$$(4.11a) \quad E \{ Y_s^* \mid N_s^* \} \leq N_s^* H(f_k(s));$$

$$(4.11b) \quad \text{there is a positive } \varepsilon \text{ (which depends on } L^* \text{ but not on } k \text{ or } n) \text{ such that } E \{ Y_s^* \mid N_s^* \} \leq N_s^* H(f_k(s)) - \varepsilon N_s^* \text{ on } \{ 2 \leq N_s^* \leq L^* \}.$$

These results follow from (3.8) in [8]. Thus,

$$E \{ Y_s^* \} \leq \lambda H(f_k(s)) - \varepsilon P \{ 2 \leq N_s^* \leq L^* \}.$$

Since $2^{-K_4} \leq \lambda \leq 2^{K_3}$, $P \{ 2 \leq N_s^* \leq L^* \} / \lambda$ is bounded above and below. There is a small positive ε' , which does not depend on k or n , such that

$$(4.12) \quad E \{ Y_s^* \} \leq \lambda [H(f_k(s)) - \varepsilon'].$$

We can now use Bernstein's inequality (3.2). Indeed, by the definition of Γ -uniformity, $\gamma_s \geq b > 0$; see (7) in [8]. And

$$-N_s^* + \log b < Y_s^* < 0,$$

by (3.3d) in [8]. Furthermore, the Y_s^* are independent. We take $m = 2^k$, $\xi = \text{Pois}(2^{K_3}) + 2 \log b + 2^{K_3}$, $y = \varepsilon'' n$, where ε'' is fixed with

$$0 < \varepsilon'' < \min \{ \varepsilon', h_1 / 2^{K_3} \},$$

so $\varepsilon'' < \varepsilon'$ and $y < h_1 m$. See (3.3) to motivate the definition of ξ . Then

$$(4.13) \quad P \left\{ \sum_{s \in C_k} Y_s^* \geq \sum_{s \in C_k} E \{ Y_s^* \} + \varepsilon'' n \right\} < r^n,$$

where $r = \exp(-Cn/m)$ and $C = \varepsilon''^2 / 2\sigma^2$. But $n/m = n/2^k \geq 2^{-K_4}$. We take $\varrho = \exp(-C2^{-K_4})$. Combine (4.12) and (4.13):

$$(4.14) \quad P \left\{ n^{-1} \sum_{s \in C_k} Y_s^* \geq \int H(f_k) - \varepsilon' + \varepsilon'' \right\} < \varrho^n.$$

We take $\delta = \varepsilon' - \varepsilon'' > 0$. This proves (4.10). ■

(4.15) COROLLARY. Suppose (1.1) and (1.4). Suppose the π_k are Γ -uniform. Fix positive integers K_3 and K_4 . Then there is a small positive δ (depending on K_3 and K_4) such that, almost surely, for all sufficiently large n , for all k with $\log_2 n - K_3 \leq k \leq \log_2 n + K_4$,

$$L_{k,n} < \int H(f_k) d\lambda^\infty - \delta.$$

This completes our discussion of the endzone.

The high zone: $\log_2 n + K_4 \leq k$.

Let

$$(4.16) \quad H(p, \theta) = p \log \theta + (1-p) \log(1-\theta).$$

The relative entropy function H is left undefined at the corners $p = \theta = 0$ or 1 , where it has singularities.

Let $s \in S_k$ iff $s \in C_k$ and $N_s = 1$: the S is for "singly-occupied." Let

$$(4.17) \quad S_{k,n} = \sum_{s \in S_k} \log \phi(N_s, X_s, \gamma_s).$$

In other words, $S_{k,n}$ represents the sum defining $L_{k,n}$, extended only over the singly-occupied s . Since $0 < \phi < 1$, we have $L_{k,n} \leq S_{k,n}/n$.

From the definition of Γ -uniformity, given as (7) in [8], g_s is the mean of γ_s ; if $k > k_1$ and $s \in C_k$, then $g_s = g_\infty(s)$, the function g_∞ being constant on each s in C_k .

(4.18) LEMMA. Suppose (1.1) and (1.4). Suppose the π_k are Γ -uniform. Then for any positive δ , there is a positive $\varrho < 1$ and a positive integer K_4 (both depending on δ) and a finite positive constant A such that, for all n and all $k \geq \log_2 n + K_4$,

$$P_f \{ S_{k,n} \geq n(e^{-\lambda} \int H(f_k, g_\infty) d\lambda^\infty + \delta) \} < A \sqrt{n} \varrho^{n^{2/2^k}}.$$

Proof. The argument is by Poissonization. Define N_s^* , X_s^* , and Y_s^* as in (4.9). It suffices to prove

$$(4.19) \quad P \left\{ \sum_{s \in C_k} \tilde{Y}_s \geq n(e^{-\lambda} \int H(f_k, g_\infty) d\lambda^\infty + \delta) \right\} < \varrho^{n^{2/2^k}},$$

where $\tilde{Y}_s = Y_s^*$ when $N_s^* = 1$, and $\tilde{Y}_s = 0$ elsewhere. Now, for $N_s^* = 1$,

$$\tilde{Y}_s = X_s^* \log g_s + (1 - X_s^*) \log(1 - g_s) = X_s^* \log g_\infty(s) + (1 - X_s^*) \log(1 - g_\infty(s)).$$

In particular,

$$E \{ \tilde{Y}_s \mid N_s^* = 1 \} = H[f_k(s), g_\infty(s)].$$

Since $P \{ N_s^* = 1 \} = \lambda e^{-\lambda}$, and $\lambda = n/2^k$, we obtain

$$(4.20) \quad E \left\{ \sum_{s \in C_k} \tilde{Y}_s \right\} = \lambda e^{-\lambda} \sum_{s \in C_k} H[f_k(s), g_\infty(s)] = n e^{-\lambda} \int H[f_k(s), g_\infty(s)].$$

The function g_∞ is bounded between b and B , with $0 < b < B < \infty$, again by definition. So the random variables \tilde{Y}_s are uniformly bounded, say by C . We use Bernstein's inequality (3.2) with $\xi \equiv 2C$, $y = n\delta$, $m = 2^k$:

$$(4.21) \quad P \left\{ \sum_{s \in C_k} \tilde{Y}_s \geq n(e^{-\lambda} \int H(f_k, g_\infty) d\lambda^\infty + \delta) \right\} < \exp \left(-\frac{n^2 \delta^2}{2\sigma^2 2^k} \right).$$

The condition $y \leq h_1 m$ is satisfied if K_4 is large enough. This proves (4.19), with $\varrho = \exp(-\delta^2/2\sigma^2)$. ■

(4.22) LEMMA. Suppose (1.1) and (1.4). Suppose the π_k are Γ -uniform, and $f \neq g_\infty$. Then there is a small positive δ and a large positive integer K_4 such that, almost surely $[P_f]$, for all sufficiently large n , for all k with $\log_2 n + K_4 \leq k$,

$$L_{k,n} < \int H(f) d\lambda^\infty - \delta.$$

Proof. $L_{k,n} \leq S_{k,n}/n$, and $S_{k,n}$ decreases with increasing k . (Eventually, $S_{k,n}$ stabilizes). The reason is that S_k , the set of singly-occupied cells, increases with k . Thus, it suffices to consider the least $k \geq \log_2 n + K_4$. We must show that almost surely, for all sufficiently large n , for the least $k \geq \log_2 n + K_4$,

$$(4.23) \quad S_{k,n}/n < \int H(f) d\lambda^\infty - \delta.$$

We choose $\delta > 0$ so small that $\int H(f, g_\infty) d\lambda^\infty < \int H(f) d\lambda^\infty - 4\delta$. Now $f_k \rightarrow f$, so for K_4 large and $k \geq \log_2 n + K_4$,

$$\int H(f_k, g_\infty) d\lambda^\infty < \int H(f, g_\infty) d\lambda^\infty + \delta.$$

But $\lambda = n/2^k \leq 1/2^{K_4}$; H is negative; for K_4 large,

$$\begin{aligned} e^{-\lambda} \int H(f_k, g_\infty) d\lambda^\infty + \delta &< e^{-\lambda} (\int H(f, g_\infty) d\lambda^\infty + \delta) + \delta \\ &< \int H(f, g_\infty) d\lambda^\infty + 3\delta < \int H(f) d\lambda^\infty - \delta. \end{aligned}$$

Now (4.18) proves (4.23), because

$$\sum_{n=1}^{\infty} \sqrt{n} \varrho^{n^{2/2^k}} < \sum_{n=1}^{\infty} \sqrt{n} \varrho^{n^{2^{K_4}}} < \infty. \quad \blacksquare$$

Discussion. For this part of the argument, we do not need that F , the set of prior means, is finite; we do need $b \leq \gamma \leq B$. We also do not need that $g_k \equiv g_\infty$ for all large k ; uniform convergence would be enough, or even convergence in measure. Finally, we do not need that g_∞ is finitary, continuous, etc.

We fix $\delta > 0$ and choose K_4 large to control the high zone, by an entropy rate argument. For any choice of K_3 and K_4 , the endzone goes away. We choose K_3 and K_2 large to control the upper midzone, in the sense of showing that $\tilde{\pi}_{k,n}$ will be close to f_k , and hence f : see (4.3)–(4.4). This may be inefficient, because the upper midzone is probably irrelevant. For any K_2 , we get consistency in the lower midzone by (4.2); and likewise for the early zone if $f = f_k$. Details are omitted because they parallel [8]. This concludes our discussion of Theorem (1.5).

The proof of Theorem (1.6). The argument is more delicate; the rate of convergence of g_k to g_∞ matters, and so does the behavior of g_∞ . For simplicity, we assume that

$$(5.1) \quad f_k = f = p \text{ for all } k \quad \text{and} \quad g_k = g_\infty = p \text{ for } k > k_1.$$

The high zone splits as follows:

Early high zone:

$$\log_2 n + K_4 \leq k \leq 2\log_2 n - \log_2 \log n - K_5.$$

Middle high zone:

$$2\log_2 n - \log_2 \log n - K_5 \leq k \leq 2\log_2 n - \log_2 \log n + K_6.$$

Late high zone:

$$2\log_2 n - \log_2 \log n + K_6 \leq k \leq 3.1\log_2 n.$$

Very late high zone:

$$3.1\log_2 n \leq k.$$

We now give some heuristics for the early zone, lower midzone, and upper midzone, that is, for $k \leq \log_2 n - K_3$:

$$(5.2a) \quad \log Q_{k,n} \doteq \sum_{s \in C_k} [N_s H(\hat{p}_s) - \frac{1}{2} \log N_s]$$

and

$$(5.2b) \quad \sum_{s \in C_k} N_s H(\hat{p}_s) \doteq nH(p) + T_n + H'(p) Q_{k,n},$$

where η_i is the response of subject i ,

$$(5.3a) \quad T_n = H'(p) \sum_{i=1}^n (\eta_i - p)$$

and

$$(5.3b) \quad Q_{k,n} = \sum_{s \in C_k} N_s (\hat{p}_s - p)^2.$$

Furthermore,

$$\sum_{s \in C_k} \log N_s \doteq 2^k \log(n/2^k).$$

(The expression $n/2^k$ represents the number of observations per parameter.) To sum up,

$$\log Q_{k,n} \doteq nH(p) + T_n - \frac{1}{2} 2^k \log(n/2^k).$$

The class of theories k with $k \leq \log_2 n - K_3$ is dominated by the smallest k with positive prior weight, namely, theory l . (In the upper midzone, another nuisance term appears in the expansion; but the argument goes through anyway.) The endzone goes away by previous arguments. The early and middle high zones can also be eliminated.

The late and very late high zones remain, and the term in $2^k \log(n/2^k)$ drops out:

$$\log Q_{k,n} \doteq nH(p) + T_n.$$

Therefore, late theories compete — on entropy grounds — with theory l . It is the rate of decay of the theory weights w_k which decides the issue. The competitive late zone starts more or less at $k = 2 \log_2 n$, when there are $1/n$ data points per parameter. In [8], the cutoff was 1 data point per parameter; the extra randomness in N_s helps the Bayesian statistician and changes the critical rate for w_k from $1/2^{k/2}$ to $1/2^{k/4}$.

Now for the details. We begin by showing that $Q_{k,n}$ is small relative to $2^k \log(n/2^k)$ provided $k \leq \log_2 n - K_3$.

(5.4) LEMMA. Define $Q_{k,n}$ by (5.3b). For each n , $Q_{k,n}$ increases with k .

Proof. Use Jensen's inequality. ■

(5.5) LEMMA. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Let K_1 be an arbitrary positive integer. Then almost surely $[P_f]$, for all sufficiently large n , for all $k \leq K_1$,

$$Q_{k,n} < 2^k \cdot 2 \cdot \log \log n.$$

Proof. Use the law of the iterated logarithm. ■

(5.6) LEMMA. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Define σ^2 and h_1 as in (3.4). Fix $B > 2$. Then there is a large positive integer K_2 (depending on B) such that: almost surely $[P_f]$, for all sufficiently large n , for all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n$,

$$Q_{k,n} < p(1-p)2^k + \sigma \sqrt{B2^k \log n}.$$

Proof. This is immediate from (3.4), with $m = 2^k$ and $y = \sigma \sqrt{B2^k \log n}$. The test sum for the Borel–Cantelli lemma is at most

$$\sum_n (\log_2 n) / n^{B/2} < \infty.$$

And the condition $y \leq h_1 m$ is satisfied if K_2 is large enough. ■

(5.7) LEMMA. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Fix $\delta > 0$. Choose K_2 as in (5.6). Then there is a large positive integer K_3 (depending on δ) such that: almost surely $[P_f]$, for all sufficiently large n , for all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$,

$$Q_{k,n} < \delta 2^k \log_2(n/2^k).$$

Proof. Suppose first that $\log_2 \log n + K_2 \leq k \leq \frac{1}{2} \log_2 n$. Write $a_n \ll b_n$ iff $a_n/b_n \rightarrow 0$. Then

$$p(1-p)2^k + \sigma \sqrt{B2^k \log n} \ll \frac{1}{2} 2^k \log_2 n \leq 2^k \log_2(n/2^k).$$

Suppose next that $\frac{1}{2}\log_2 n \leq k \leq \log_2 n - K_3$. Then

$$p(1-p)2^k + \sigma \sqrt{B2^k \log n} \leq \delta 2^k \log_2(n/2^k)$$

provided K_3 is large. Indeed, $p(1-p) \leq 1/4$ and $\log_2(n/2^k) \geq K_3$ which is large, taking care of the term $p(1-p)2^k$. Finally, $\sigma \sqrt{B2^k \log n} \leq 2^k$. ■

(5.8) LEMMA. Assume (1.1), (1.4) and (5.1). Suppose the π_k are Γ -uniform. Fix $\delta > 0$. Choose K_2 as in (5.6). Then there is a large positive integer K_1 (depending on δ) such that: almost surely $[P_f]$, for all sufficiently large n , for all k with $K_1 \leq k \leq \log_2 \log n + K_2$,

$$Q_{k,n} < \delta 2^k \log_2(n/2^k).$$

Proof. Let k_n be the least positive integer which is $\log_2 \log n + K_2$ or more. Now

$$\begin{aligned} Q_{k,n} &\leq Q_{k_n,n} && \text{by (5.4)} \\ &< p(1-p)2^{k_n} + \sigma \sqrt{B2^{k_n} \log n} && \text{by (5.6)} \\ &\leq (p(1-p)2^{K_2+1} + \sigma \sqrt{B2^{K_2+1}}) \log n < \delta 2^k \log_2(n/2^k) \end{aligned}$$

for k with $K_1 \leq k \leq \log_2 \log n + K_2$ provided K_1 is large. (The 1st and 3rd inequalities hold for all n ; the 2nd and 4th for n large.) ■

(5.9) COROLLARY. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Fix $\delta > 0$. Choose K_3 as in (5.7). Then almost surely $[P_f]$, for all sufficiently large n , for all k with $k \leq \log_2 n - K_3$,

$$Q_{k,n} < \delta 2^k \log_2(n/2^k).$$

Note. From here on, K_3 is forced large; but K_1 and K_2 are free again.

Proof. Combine (5.5), (5.7), and (5.8). ■

This completes the discussion of $Q_{k,n}$, and we turn to the term $\sum_{s \in C_k} \log N_s$ in the expansion (5.2a) of $\log Q_{k,n}$. The sum is $[1 + o(1)] 2^k \log \lambda$, where $\lambda = n/2^k$ as in (4.1). The main technique is Poissonization, to approximate the ball-dropping distribution (2.3a). Unfortunately, there are zones which do not quite match those previously defined. We begin with $k \leq (\log_2 n)/4$.

(5.10) LEMMA. Assume (1.4). Fix $\delta > 0$. Then for all n , all $k \leq (\log_2 n)/2$, and all $s \in C_k$,

- (a) $P_f \{N_s/\lambda \geq 1 + \delta\} < \exp(-\delta^2 \sqrt{n}/2)$,
- (b) $P_f \{N_s/\lambda \leq 1 - \delta\} < \exp(-\delta^2 \sqrt{n}/2)$.

Proof. As (2.3a) shows, N_s is $\text{bin}(n, 1/2^k)$. Now use (3.1). Of course, $\lambda^2/2^k = n^2/2^{3k} \geq \sqrt{n}$ since $k \leq (\log_2 n)/2$. ■

(5.11) LEMMA. Assume (1.4). Fix $\delta > 0$. Then almost surely $[P_f]$, for all sufficiently large n and all $k \leq (\log_2 n)/2$,

$$\left| \sum_{s \in C_k} (\log N_s - \log \lambda) \right| < \delta 2^k \log \lambda.$$

Proof. By (5.10) and the Borel–Cantelli lemma, $1 - \delta \leq N_s/\lambda \leq 1 + \delta$ for all $s \in C_k$ and all $k \leq (\log_2 n)/2$, for all sufficiently large n , almost surely: the test sum is bounded by

$$2 \sum_n \sum_{k=0}^{(\log_2 n)/2} 2^k \exp(-\delta^2 \sqrt{n}/2) < 4 \sum_n \sqrt{n} \exp(-\delta^2 \sqrt{n}/2) < \infty.$$

Finally, $k \leq (\log_2 n)/2$ entails

$$2^k |\log(1 \pm \delta)| \leq 2^k \log(n/2^k). \quad \blacksquare$$

We turn now to larger k ; the lower endpoint of the range is not material, but $\log_2 \log n$ is a convenient cut-point.

(5.12) LEMMA. Assume (1.4). For $s \in C_k$, let N_s^* be independent $\text{Pois}(\lambda)$ variables. Let $\tilde{N}_s = N_s^*$ when $N_s^* > \lambda e^{-1/2}$, else let $\tilde{N}_s = \lambda e^{-1/2}$.

(a) Fix $B > 2$. Then there is a positive integer K_2 so large (depending on B) that, for all n and all $k \geq \log_2 \log n + K_2$,

$$P \left\{ \left| \sum_{s \in C_k} (\log \tilde{N}_s - E \{ \log \tilde{N}_s \}) \right| \geq B \sigma 2^k \sqrt{\frac{\log n}{n}} \right\} < 2/n^{B^2/2}.$$

(b) Fix $\delta > 0$ and $C > 2$. Then there are positive integers K_2 and K_3 so large (depending on δ and C) that, for all n and all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$, the chance that $N_s^* \leq \lambda e^{-1/2}$ for $\delta 2^k$ or more indices $s \in C_k$ is bounded above by $1/n^C$.

Proof. Claim (a). This follows from (3.11) with $m = 2^k$, all $\lambda_i = \lambda = n/2^k$, and $y = B \sigma \sqrt{2^k \log n}$. The condition $y \leq h_1 m$ is satisfied if K_2 is large.

Claim (b). By (3.5b), $P \{ N_s^* \leq \lambda e^{-1/2} \} < e^{-\lambda/16} < \delta/2$ provided K_3 is large; indeed, $\lambda \geq 2^{K_3}$. The chance that $\delta 2^k$ or more of these unlikely events occur can be bounded above by (3.1b). The bound is

$$\exp(-\delta^2 2^k/8) \leq \exp(-\delta^2 2^{K_2-3} \log n)$$

because $2^k/8 \geq 2^{K_2-3} \log n$. \blacksquare

(5.13) LEMMA. Assume (1.4). Let $N'_s = N_s$ when $N_s > \lambda e^{-1/2}$, else let $N'_s = \lambda e^{-1/2}$. Fix $\delta > 0$, choose K_2 and K_3 as in (5.12). Then:

(a) Almost surely $[P_f]$, for all sufficiently large n , and all k with

$$\log_2 \log n + K_2 \leq k \leq \log_2 n - 1,$$

we have

$$\left| \sum_{s \in C_k} (\log N'_s - \log \lambda) \right| < \delta 2^k \log \lambda.$$

(b) Almost surely $[P_f]$, for all sufficiently large n , and all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$,

$$0 \leq \sum_{s \in C_k} (\log N'_s - \log N_s) < \delta 2^k \log \lambda.$$

Proof. Claim (a). We de-Poissonize (5.12a):

$$(5.14) \quad P_f \left\{ \left| \sum_{s \in C_k} (\log N'_s - E \{ \log \tilde{N}_{s'} \}) \right| \geq B \sigma 2^k \sqrt{\frac{\log n}{n}} \right\} < A/n^{(B^2-1)/2}.$$

By (5.14) and the Borel-Cantelli lemma, almost surely, for all sufficiently large n , for all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - 1$,

$$\left| \sum_{s \in C_k} (\log N'_s - E \{ \log \tilde{N}_{s'} \}) \right| < B \sigma 2^k \sqrt{\frac{\log n}{n}};$$

indeed, the test sum is bounded above by

$$A \sum_n (\log_2 n) / n^{(B^2-1)/2} < \infty.$$

Since $k \leq \log_2 n - 1$, we have

$$B \sigma 2^k \sqrt{(\log n)/n} \leq 2^k \log_2 (n/2^k).$$

Now use (3.12).

Claim (b). We de-Poissonize (5.12b). Let $s \in S_k$ iff $s \in C_k$ and $N_s \leq \lambda e^{-1/2}$: the S is for "small." Write $|S_k|$ for the cardinality of S_k . Now $P_f \{ |S_k| \geq \delta 2^k \} < A/n^{C-0.5}$. There are at most $\log_2 n$ indices k to consider, and $\sum (\log_2 n) n^{C-0.5} < \infty$ because $C > 2$. Thus, almost surely, for all sufficiently large n , for all k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$, $|S_k| < \delta 2^k$.

If $s \notin S_k$, then $N'_s = N_s$. Now suppose $s \in S_k$. If $N_s = 0$, then $\log N_s = 0$ by definition. Thus,

$$0 \leq \log N'_s - \log N_s \leq \log \lambda - 1/2 < \log \lambda.$$

Consequently,

$$0 \leq \sum_{s \in C_k} (\log N'_s - \log N_s) < |Q_k| \log \lambda < \delta 2^k \log \lambda. \quad \blacksquare$$

(5.15) Remark. Assume (1.4). Fix $L > 6$. Then almost surely, for all sufficiently large n , for all k with $k \leq \log_2 n - \log_2 \log n - L$, and all $s \in C_k$, $N'_s = N_s$.

Proof. This follows from (3.5b) and Poissonization:

$$P \{ N_s \leq \lambda e^{-1/2} \} \leq A \sqrt{n} e^{-\lambda/16} \leq A/n^C, \quad \text{where } C > 2.$$

The test sum for the Borel-Cantelli lemma is bounded above by

$$A \sum_n \sum_{k=0}^{\log_2 n} 2^k / n^C < 2A \sum_n 1/n^{C-1} < \infty. \quad \blacksquare$$

(5.16) COROLLARY. Assume (1.4). Fix $\delta > 0$. Then almost surely $[P_f]$, for all sufficiently large n , for all k with $k \leq \log_2 n - K_3$,

$$\left| \sum_{s \in C_k} (\log N_s - \log \lambda) \right| < \delta 2^k \log \lambda.$$

Proof. For k with $\log_2 \log n + K_2 \leq k \leq \log_2 n - K_3$, use (5.13). For $k \leq \log_2 \log n + K_2$, use (5.11). ■

In the early zone and lower midzone, $k \leq \log_2 \log n + K$; then \hat{p}_s is nearly p : see (4.2). In these zones, we can estimate $\log \varrho_{k,n}$ as follows.

(5.17) PROPOSITION. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Define T_n by (5.3a). Fix $\delta > 0$ and $K < \infty$. Then almost surely, for all sufficiently large n , for all k with $0 \leq k \leq \log_2 \log n + K$,

$$|\log \varrho_{k,n} - nH(p) - T_n + \frac{1}{2} 2^k \log(n/2^k)| < \delta 2^k \log(n/2^k).$$

Proof. We estimate $\log \varrho_{k,n}$ by (3.3) in [8], making (5.2a) rigorous by adding $O(2^k) = o(2^k \log(n/2^k))$. Now

$$\sum_{s \in C_k} N_s H(\hat{p}_s)$$

can be expanded around p by (3.14) in [8]. The lead term is $nH(p)$. The linear term gives T_n after a bit of algebra. The quadratic remainder is negligible by (5.9). Finally,

$$\frac{1}{2} \sum_{s \in C_k} \log N_s$$

can be estimated by (5.16). ■

In the upper midzone, N_s may be 0 for some s . The corresponding terms contribute 0 to the sum defining $\log \varrho_{k,n}$. Even if $N_s > 0$, \hat{p}_s may be 0 or 1. This necessitates some additional nuisance terms in the expansion of $\log \varrho_{k,n}$ because the approximation to $\phi(m, j, \gamma)$ changes when $j = 0$ or m . See (3.2) and (5.4) in [8].

(5.18a) Let N_k be the number of $s \in C_k$ with $N_s > 0$ and $X_s = 0$ or N_s .

(5.18b) Let $s \in G_k$ iff $0 < X_s < N_s$.

(5.18c) Let $\mathcal{E}_{k,n} = -\frac{1}{2} \log(n/2^k) N_k + \sum_{s \in G_k} \log \sqrt{\hat{p}_s (1 - \hat{p}_s)}$.

All terms in $\mathcal{E}_{k,n}$ are negative because $0 < \hat{p}_s < 1$.

(5.19) PROPOSITION. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Fix $\delta > 0$ and $K < \infty$. Define K_3 as in (5.16). Then almost surely $[P_f]$, for all sufficiently large n , for all k with $\log_2 \log n + K \leq k \leq \log_2 n - K_3$,

$$|\log \varrho_{k,n} - nH(p) - T_n + \frac{1}{2} 2^k \log(n/2^k) - \mathcal{E}_{k,n}| < \delta 2^k \log(n/2^k).$$

Proof. This is argued like (5.17). ■

This completes the discussion of the early zone and midzone. The endzone goes away by (4.15), and we turn to the high zone.

The early high zone. The early high zone is defined by the condition

$$(5.20) \quad \log_2 n + K_4 \leq k \leq 2\log_2 n - \log_2 \log n - K_5.$$

K_4 defines the right edge of the endzone, but from our perspective, it is a free parameter: (4.15) imposed no condition on K_4 . For present purposes too, K_4 is not really material; we can set $K_4 = 3$. We will prove:

(5.21) PROPOSITION. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Fix a large positive number L . Then there is a large positive integer K_5 such that: almost surely $[P_f]$, for all sufficiently large n , for all k satisfying (5.20),

$$\log Q_{k,n} < nH(p) + T_n - L \log n.$$

Suppose $s \in C_k$. As in (4.17), let $s \in S_{k,n}$ iff $N_s = 1$; likewise, $s \in M_{k,n}$ iff $N_s > 1$. The S is for "single occupancy," and M for "multiple occupancy"; the dependence on n will matter later. Write $i \in s$ iff $\xi_j(i) = s_j$ for $1 \leq j \leq k$; in other words, the first k covariates for subject i agree with s . Suppose k is so large that $g_k \equiv p$: see (5.1). A bit of algebra shows

$$(5.22) \quad \text{If } s \in S_{k,n}, \text{ then } \log \phi(N_s, X_s, \gamma_s) = H(p) + (\eta_i - p)H'(p) \text{ for the unique } i \in s.$$

For $0 \leq j \leq m$ and $m \geq 2$, let

$$(5.23) \quad \phi_0(m, j, \gamma) = \log \phi(m, j, \gamma) - mH(p) - (j - mp)H'(p).$$

For $s \in M_{k,n}$, let $\Delta_s = \phi_0(N_s, X_s, \gamma)$. By (5.22) and a bit more algebra,

$$(5.24) \quad \log Q_{k,n} = nH(p) + T_n + \sum_{s \in M_{k,n}} \Delta_s.$$

To prove (5.21), we must estimate $\sum_{s \in M_{k,n}} \Delta_s$. The main technique is Poissonization, and here are some preliminaries. Recall from the definition (7) in [8] of Γ -uniformity that $\gamma \in \Gamma$ entails $\gamma \geq b > 0$. The next result is immediate from (3.3d) in [8].

$$(5.25) \text{ LEMMA. } |\phi_0(m, j, \gamma)| \leq [1 + |H'(p)|]m + |\log b| \text{ for } \gamma \in \Gamma \text{ with lower bound } b.$$

(5.26) DEFINITION. Fix k . For $s \in C_k$, let N_s^* be $\text{Pois}(\lambda)$, where $\lambda = n/2^k$. Given $\{N_s^*\}$, let $\{X_s^*\}$ be independent $\text{bin}(N_s^*, p)$. Let $\Delta_s^* = \phi_0(N_s^*, X_s^*, \gamma_s)$. Let M^* be the number of $s \in C_k$ with $N_s^* \geq 2$.

Relationships (5.27)–(5.30) are obvious:

$$(5.27) \quad M^* = \sum_{s \in C_k} I_s^*, \text{ where } I_s^* \text{ is } 0 \text{ if } N_s^* < 2 \text{ and } I_s^* \text{ is } 1 \text{ if } N_s^* \geq 2. \text{ The } I_s^* \text{ are iid.}$$

(5.28) $\frac{1}{2}\lambda^2(1-\lambda) \leq [1-(1+\lambda)e^{-\lambda}] \leq \frac{1}{2}\lambda^2$ for all λ .

(5.29) $E\{M^*\} = 2^k [1-(1+\lambda)e^{-\lambda}]$.

(5.30) $\frac{1}{2}n\lambda(1-\lambda) \leq E\{M^*\} \leq \frac{1}{2}n\lambda$ for all λ .

(5.31) LEMMA. Fix δ with $0 < \delta < 1$. Suppose $0 < \lambda < \delta/2$. Then:

(a) $P\{M^* \geq (1+\delta)n\lambda/2\} < \exp(-\delta^2n\lambda/8)$;

(b) $P\{M^* \leq (1-\delta)n\lambda/2\} < \exp(-\delta^2n\lambda/16)$.

Proof. Claim (a). This is Bernstein's inequality. Theorem (4) in [13], coupled with the estimate (5.30) for $E\{M^*\}$, gives the bound

$$\exp\left(-\frac{1}{2} \frac{(\delta n\lambda/2)^2}{(1+\delta)n\lambda/2}\right) < \exp(-\delta^2n\lambda/8)$$

because $0 < \delta < 1$.

Claim (b) is similar. By (5.30),

$$n\lambda/2 \geq E\{M^*\} \quad \text{and} \quad E\{M^*\} - (1-\delta)n\lambda/2 \geq n\lambda(\delta-\lambda)/2 > \delta n\lambda/4,$$

so the bound is

$$\exp\left(-\frac{1}{2} \frac{(\delta n\lambda/4)^2}{n\lambda/2}\right) < \exp(-\delta^2n\lambda/16). \quad \blacksquare$$

Note. Lemma 3.1 is quite inefficient for small p , when \sqrt{mx} would — ideally — be replaced by \sqrt{mpx} . Hence the resort to other estimates.

(5.32) LEMMA. For $i = 1, 2, \dots$ let \tilde{A}_i be independent and distributed as Δ_s^* , given $N_s^* \geq 2$. Define $\varepsilon_m > 0$ as in (3.8) of [8]. Suppose $0 < \lambda < 1/2$. Then:

(a) $E\{\tilde{A}_i\} < -\varepsilon_2$.

(b) There is an $\varepsilon > 0$ and σ^2 with $0 < \sigma^2 < \infty$ such that:

(i) for all λ with $0 < \lambda < 1/2$ and all $m = 1, 2, \dots$

$$P\left\{\sum_{i=1}^m \tilde{A}_i \geq -\varepsilon m\right\} < \exp(-\varepsilon^2 m/2\sigma^2);$$

(ii) $P\left\{\sum_{s \in C_k} \Delta_s^* I_s^* \geq -\varepsilon M^*\right\} < \exp(-\varepsilon^2 n\lambda/8\sigma^2)$.

Proof. Claim (a). By (3.8) in [8], $E\{\Delta_s^* \mid N_s^* = m\} < -m\varepsilon_m$ for $m \geq 2$. So

$$\begin{aligned} E\{\tilde{A}_i\} &= E\{\Delta_s^* \mid N_s^* \geq 2\} < -2\varepsilon_2 P\{N_s^* = 2\} / P\{N_s^* \geq 2\} \\ &< -2\varepsilon_2(1-\lambda) < -\varepsilon_2, \end{aligned}$$

with the help of (5.28).

Claim (b). Let $\xi = |\log b| + \{1 + |H(p)| + |H'(p)|\} N$, where N is Poisson(1/2) conditioned to be 2 or more. By (5.25) and (3.25), $|\tilde{A}_i|$ is stochastically bounded by ξ , so Lemma (3.2) applies. Compute σ^2 and h_1 according to that lemma.

Let $\varepsilon = \min \{h_1, \varepsilon_2/2\}$ and $\varrho = \exp(-\varepsilon^2/2\sigma^2)$. Now

$$P \left\{ \sum_{i=1}^m \tilde{\Delta}_i \geq -\varepsilon m \right\} < P \left\{ \sum_{i=1}^m \tilde{\Delta}_i \geq \sum_{i=1}^m E \{ \tilde{\Delta}_i \} + \varepsilon m \right\} < \exp(-\varepsilon^2 m/2\sigma^2).$$

The first inequality holds because $E \{ \tilde{\Delta}_i \} < -2\varepsilon$; the second, by (3.2): the condition $y_1 \leq h_1 m$ holds because $\varepsilon \leq h_1$. This proves (i), and we turn to (ii).

Conditional on $M^* = m$, the sum is distributed as $\sum_{i=1}^m \tilde{\Delta}_i$, giving the bound

$$E \{ \exp(-\varepsilon^2 M^*/2\sigma^2) \} < \exp[-\varepsilon^2 E \{ M^* \} / 2\sigma^2]$$

by Jensen's inequality. But $\lambda < 1/2$, so $E \{ M^* \} > n\lambda/4$ by (5.30). ■

Note. In the proof of (b), if you just think of $\sum_{s \in C_k} \Delta_s^* I_s^*$ as the sum of 2^k terms, (3.2) gives the disappointing bound

$$\exp \left(-\frac{\varepsilon^2 (n\lambda)^2}{8\sigma^2 2^k} \right) = \exp(-\varepsilon^2 n\lambda^3/8\sigma^2).$$

Recall that $M_{k,n} = \{s: s \in C_k \text{ and } N_s \geq 2\}$.

(5.33) LEMMA. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Fix δ with $0 < \delta < 1/2$ and suppose $0 < \lambda < \delta/2$. Choose $\varepsilon > 0$ as in (5.32b). Then:

- (a) $P_f \{ |M_{k,n}| \geq (1 + \delta)n\lambda/2 \} < A\sqrt{n} \exp(-\delta^2 n\lambda/8)$;
- (b) $P_f \{ |M_{k,n}| \leq (1 - \delta)n\lambda/2 \} < A\sqrt{n} \exp(-\delta^2 n\lambda/16)$;
- (c) $P_f \{ \sum_{s \in M_{k,n}} \Delta_s \geq -\varepsilon |M_{k,n}| \} < A\sqrt{n} \exp(-\varepsilon^2 n\lambda/8\sigma^2)$.

Proof. Claims (a) and (b) follow from (5.31) by de-Poissonization. Claim (c) is similar, starting from (5.32). ■

(5.34) COROLLARY. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Suppose k and n satisfy (5.20). Then:

- (a) $P_f \{ |M_{k,n}| \leq (1 - \delta)n\lambda/2 \} < A/n^{C-0.5}$ with $C = \delta^2 2^{K_5-4}$;
- (b) $P_f \{ \sum_{s \in M_{k,n}} \Delta_s \geq -\varepsilon |M_{k,n}| \} < A/n^{D-0.5}$ with $D = \varepsilon^2 2^{K_5-3}/\sigma^2$.

Proof of Proposition (5.21). Fix $\delta < 1/2$ in (5.34a); we require $2^{-K_4} < \delta/2$, so $\lambda < \delta/2$ in (5.31). Choose ε as in (5.32b). Choose K_5 so large that $C > 2$ and $D > 2$ in (5.34). There are at most $\log_2 n$ theories in the zone. So, almost surely, for all sufficiently large n , for all k satisfying (5.20),

$$|M_{k,n}| > (1 - \delta)n\lambda/2 > 2^{K_5-2} \log n,$$

$$\sum_{s \in M_{k,n}} \Delta_s < -\varepsilon |M_{k,n}| < -\varepsilon 2^{K_5-2} \log n. \quad \blacksquare$$

Remark. We have assumed in (5.1) that the mean of γ_s equals p for $s \in C_k$ and $k > n_1$. Suppose that γ_s is constant, say at $\gamma \in \Gamma$ with $\int \theta\gamma(\theta) d\theta = p$. Then $\{ \Delta_s^*: s \in M_{k,n} \}$ are iid for each k , with $E \{ \Delta_s^* \} = E \{ \phi_0(N_\lambda, X, \gamma) \}$; N_λ is $\text{Pois}(\lambda)$ conditioned to be 2 or more; given $N_\lambda = m$, X is $\text{bin}(m, p)$; see (5.23) and (5.26).

The argument for (5.33) shows that

$$\sum \{A_s: s \in M_{k,n}\} \approx \frac{1}{2} n \lambda E \{ \phi_0(N_\lambda, X, \gamma) \}.$$

This completes our discussion of the early high zone.

The middle high zone. The middle high zone is the most delicate of all the zones. It is defined by the condition

$$(5.35) \quad 2 \log_2 n - \log_2 \log n - K_5 \leq k \leq 2 \log_2 n - \log_2 \log n + K_6,$$

where K_5 and K_6 are large positive integers: K_5 is needed to control the early high zone, and K_6 will control the late high zone.

(5.36) PROPOSITION. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Then there is a small positive ε_0 (depending on K_5 and K_6) such that: almost surely, for all sufficiently large n , for all k satisfying (5.35),

$$\log Q_{k,n} < nH(p) + T_n - \varepsilon_0 \log n.$$

At stage n of the trial, we have data on n subjects; let $D_{k,n}$ be the set of $s \in C_k$ with $N_s = 2$; the D is for "doubly occupied." The main difficulty is showing that $|D_{k,n}| \approx n\lambda/2$. The dependence on n matters, and is displayed in the notation. Since $n\lambda$ is of order $\log n$, exponential bounds must be supplemented by passing to geometric subsequences, and the \sqrt{n} for de-Poissonization cannot be afforded. We solve the latter problem first.

(5.37) LEMMA. At stage n , let $D_{k,n}$ be the set of $s \in C_k$ with $N_s = 2$, and let $M_{k,n}$ be the set of $s \in C_k$ with $N_s \geq 2$. Fix δ with $0 < \delta < 1$. Fix positive integers K_5 and K_6 . Then almost surely, for all sufficiently large n , for all k satisfying (5.35),

$$(a) \quad (1 - \delta) n\lambda/2 < |M_{k,n}| < (1 + \delta) n\lambda/2,$$

$$(b) \quad (1 - \delta) n\lambda/2 < |D_{k,n}| < (1 + \delta) n\lambda/2.$$

Note. λ depends on k and n , and $n\lambda = n^2/2^k$; see (4.1).

Proof. Claim (a). Fix r slightly bigger than 1 and consider the sequence r^j . For each k , $|M_{k,n}|$ increases with n . As n increases from r^j to r^{j+1} , $n\lambda$ increases from $r^{2j}/2^k$ to $r^{2j+2}/2^k$, i.e., only by a factor of r^2 . Thus, it suffices to prove claim (a) for n of the form r^j . Recall S_n from (3.22). By (3.20), $|M_{k,n}| = S_n$ or $S_n - 1$; and it is enough to prove the claim for S_n and n of the form r^j . That is immediate from the Borel-Cantelli lemma and (3.24) with $n = r^j$ and $b = 2^k$.

Claim (b) follows from (a), because $|D_{k,n}| = |M_{k,n}|$ or $|M_{k,n}| - 1$ by (3.20).

Recall the function $\phi_0(m, j, \gamma)$ from (5.23). Let \mathcal{X} be the class of random variables distributed as $\phi_0(2, X, \gamma)$, where X is $\text{bin}(2, p)$ and $\gamma \in \Gamma$. If $Y \in \mathcal{X}$, then Y is uniformly bounded by (5.25), and $E\{Y\} < -2\varepsilon_2 < 0$ by (3.8) in [8]. As (3.2) shows,

(5.38) LEMMA. *There are positive constants h_1 and σ^2 , depending only on \mathcal{K} , such that: if $Y_i \in \mathcal{K}$ are independent for $i = 1, \dots, m$ and $0 < y \leq h_1 m$, then*

$$P \left\{ \sum_{i=1}^m (Y_i - E \{Y_i\}) \geq y \right\} < \exp(-y^2/2\sigma^2).$$

Recall Δ_s , as defined for (5.24). Given $D_{k,n}$, $\{\Delta_s: s \in D_{k,n}\}$ are independent; and $\Delta_s \in \mathcal{K}$.

(5.39) LEMMA. *Define ε_2 as in (3.8) of [8]; h_1 and σ^2 as in (5.38). Let $0 < \varepsilon < \min(h_1, \varepsilon_2)$. Then almost surely, for all sufficiently large n , for all n and k satisfying (5.35),*

$$\sum \{ \Delta_s: s \in D_{k,n} \} < -\varepsilon(\log_2 n)/2^{K_6+4}.$$

Proof. First, consider only n of the form 2^j . By (5.38),

$$(5.40) \quad P_f \left\{ \sum \{ \Delta_s: s \in D_{k,n} \} \geq -\varepsilon |D_{k,n}| \mid D_{k,n} \right\} < \exp(-\varepsilon^2 |D_{k,n}|/2\sigma^2).$$

Then

$$E \{ \exp(-\varepsilon^2 |D_{k,n}|/2\sigma^2) \} < \exp(-\varepsilon^2 E \{ |D_{k,n}| \} / 2\sigma^2)$$

and

$$E \{ |D_{k,n}| \} \approx 2^k \lambda^2 / 2 = n^2 / 2^{k+1} > (\log n) / 2^{K_6+1} > j / 2^{K_6+2}.$$

The Borel-Cantelli lemma shows that almost surely, for all sufficiently large n of the form 2^j , for all k satisfying (5.35), $\sum \{ \Delta_s: s \in D_{k,n} \} < -\varepsilon |D_{k,n}|$. Indeed, the test sum is bounded by

$$(1 + K_5 + K_6) \sum_j \exp(-C_0 j) < \infty, \quad \text{where } C_0 = \varepsilon^2 / 2^{K_6+3} \sigma^2.$$

Finally, use (5.37) to bound $|D_{k,n}|$, noting that $n\lambda \geq (\log n) 2^{K_6}$.

We must now interpolate between 2^j and 2^{j+1} ; the argument is only sketched. Fix k . Then $\sum \{ \Delta_s: s \in D_{k,n} \}$ is below $-\varepsilon(\log n) / 2^{K_6+2}$, say, when $n = 2^j$. As n increases from 2^j to 2^{j+1} , 2^j additional balls are dropped at random into the 2^k boxes, perturbing the sum. We claim that almost surely, for all sufficiently large n , the perturbations will amount at most to $\varepsilon(\log n) / 2^{K_6+3}$, so we have

$$\sum \{ \Delta_s: s \in D_{k,n} \} < -\varepsilon(\log n) / 2^{K_6+3}$$

for all n and k satisfying (5.35), with $2^j \leq n \leq 2^{j+1}$, provided j is sufficiently large.

There are four kinds of perturbations:

- (i) an additional doubly-occupied box is created, adding an independent term $\Delta_i \in \mathcal{K}$;
- (ii) a triply-occupied box may be created;
- (iii) more than one triply-occupied boxes may be created;
- (iv) a box may become more than triply occupied.

Perturbations (iii) and (iv) do not occur for large n , by (3.20), and need not be considered further. Perturbation (ii) changes the sum by a uniformly bounded amount; see (5.25).

We must now bound the effect of perturbations of type (i), showing they amount to less than $C_0 \log_2 n = C_0 j$, where $C_0 = \varepsilon/2^{K_6+5}$; this leaves more than enough to absorb perturbations of type (ii). Now, dropping in 2^j balls increases $D_{k,n}$ from (essentially) $2^{2j}/2^{k+1}$ to $2^{2j+2}/2^{k+1}$, by (5.37); i.e., from cj to $4cj$. But, adding this number of Δ 's — or any other — crosses the $C_0 j$ boundary with probability at most ρ^j , by (3.16). ■

Proof of Proposition (5.36). Use (5.24), (3.20), and (5.39). ■

Remark. We have assumed in (5.1) that the mean of γ_s equals p for $s \in C_k$ and $k > n_1$. Suppose that γ_s is constant, say at $\gamma \in \Gamma$ with $\int \theta \gamma(\theta) d\theta = p$. Then $\{\Delta_s: s \in D_{k,n}\}$ are iid with $E\{\Delta_s\} = E\{\phi_0(2, X, \gamma)\}$, X being bin(2, p); see (5.23). The argument for (5.39), pushed a little harder, shows that

$$\sum \{\Delta_s: s \in D_{k,n}\} \approx \frac{1}{2} n \lambda E\{\phi_0(2, X, \gamma)\};$$

the idea is to split along the geometric sequence r^n with r just bigger than 1.

This completes our discussion of the middle high zone.

The late high zone. The late high zone is defined by the condition

$$(5.41) \quad 2\log_2 n - \log_2 \log n + K_6 \leq k \leq 3.1\log_2 n.$$

(5.42) **PROPOSITION.** Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Fix $\varepsilon > 0$. Then there is a large positive integer K_6 (depending on ε) such that: almost surely, for all sufficiently large n , for all k satisfying (5.41),

$$|\log \rho_{k,n} - nH(p) - T_n| < \varepsilon \log n.$$

Proof. By (5.24), it is enough to bound $\sum \{\Delta_s: s \in M_{k,n}\}$. But (3.20) shows that $M_{k,n}$ (the set of multiply-occupied cells) differs from $D_{k,n}$ (the set of doubly-occupied cells) by at most one triply-occupied cell. So it is enough to bound $\sum \{\Delta_s: s \in D_{k,n}\}$, and hence $|D_{k,n}|$ by (5.25). For each n , $|D_{k,n}|$ decreases as k increases, so it is enough to consider k just larger than $2\log_2 n - \log_2 \log n + K_6$. Now (5.37) shows that almost surely, for all sufficiently large n ,

$$|D_{k,n}| < n\lambda = n^2/2^k \leq 2^{-K_6} \log n. \quad \blacksquare$$

The very late high zone. The very late high zone is defined by the condition

$$(5.43) \quad 3.1\log_2 n \leq k.$$

(5.44) PROPOSITION. Assume (1.1), (1.4), and (5.1). Suppose the π_k are Γ -uniform. Then almost surely, for all sufficiently large n , for all $k \geq 3.1 \log_2 n$,

$$\log Q_{k,n} = nH(p) + T_n.$$

Proof. This is immediate from (3.21) and (5.24). ■

Proof of Theorem (1.6). Claim (a). By (5.17), with δ_1 any small positive number of our choice,

$$(5.45) \quad \tilde{w}_{l,n} > w_l \exp [nH(p) + T_n - \frac{1}{2} 2^l \log(n/2^l) - \delta_1 2^l \log(n/2^l)].$$

The random term T_n was defined by (5.3a), and is of order $\sqrt{n \log \log n}$ or less. We must now eliminate theories in the endzone and high zone.

Theories in the endzone ($\log_2 n - K_3$ to $\log_2 n + K_4$) are negligible relative to theory l , by (4.15). The δ there depends on K_3 and K_4 ; but no matter what that δ is, the endzone has entropy rate $H(p) - \delta < H(p)$.

The early high zone is defined by (5.20). For such theories, by (5.21),

$$\begin{aligned} \sum_k \tilde{w}_{k,n} &< \left(\sum_{k=\log_2 n}^{\infty} w_k \right) \exp(nH(p) + T_n - L \log n) \\ &< \exp \left(nH(p) + T_n - L \log n - \frac{1}{4} 2^l \log n - \frac{\delta_0}{\log 2} 2^l \log n \right) \end{aligned}$$

by using the condition of the theorem — with $\log_2 n = (\log n)/\log 2$ in place of n . (The sum of the high-zone weights starts at $\log_2 n$.) In total, the early high zone has negligible posterior weight, relative to theory l , provided

$$L + \frac{\delta_0}{\log 2} 2^l > \frac{1}{4} 2^l + \delta_1 2^l.$$

But L can be made large by choosing K_5 large.

We combine the middle high zone, late high zone and very late high zone, i.e., we consider all theories

$$k \geq L(n) = 2 \log_2 n - \log_2 \log n - K_5.$$

The posterior weight in this combined zone is by (5.36), (5.42), and (5.44) at most

$$(5.46) \quad \begin{aligned} \sum_k \tilde{w}_{k,n} &< \left(\sum_{k=L(n)}^{\infty} w_k \right) \exp(nH(p) + T_n + \varepsilon \log n) \\ &< \exp(nH(p) + T_n + \varepsilon \log n - \frac{1}{4} (\log 2) 2^l L(n) - \delta_0 2^l L(n)). \end{aligned}$$

Now

$$(5.47) \quad \frac{1}{4} (\log 2) 2^l L(n) = \frac{1}{2} 2^l \log n - C_0 \log \log n - C_1.$$

Again, this zone is negligible relative to theory l if we choose

$$\varepsilon + \delta_1 2^l < \frac{2\delta_0}{\log 2} 2^l.$$

The δ_1 in (5.45) is the δ of (5.17), and is at our disposition. The ε in (5.46) comes from (5.42). To make it small, we have to choose K_6 large. Choosing K_5 and K_6 large makes the ε_0 in (5.36) small. However, the value of ε_0 does not matter.

The balance of the argument for claim (a) is omitted, being very similar to the argument for Theorem (1.5) in this paper, or Theorems 8 and 9 in [8]. Basically, posterior mass shifts into the early zone or midzone, where there are lots of observations per parameter.

Claim (b). Consider only n with

$$\sum_{k=n}^{\infty} w_k > \exp(-\frac{1}{4}(\log 2)n2^l + \delta_0 n2^l).$$

We combine theories in the late and very late high zones, so

$$k \geq L(n) = 2\log_2 n - \log_2 \log n + K_6.$$

By (5.42) and (5.44), the total posterior weight in these two zones is at least

$$(5.48) \quad \exp(nH(p) + T_n - \varepsilon \log n - \frac{1}{4}(\log 2)2^l L(n) + \delta_0 2^l L(n)).$$

ε is a small positive number, at our disposition; δ_0 is fixed by the conditions of the theorem. Of course,

$$(5.49) \quad \frac{1}{4}(\log 2)2^l L(n) = \frac{1}{2}2^l \log n - C_0 \log \log n + C_1.$$

By comparison, the total posterior weight in the early zone and midzone ($k \leq \log_2 n - K_3$) is by (5.17) and (5.19) at most

$$(5.50) \quad \sum_{k=l}^{\log_2 n - K_3} (w_k \cdot \exp[nH(p) + T_n - \frac{1}{2}2^k \log(n/2^k) + \delta_1 2^k \log(n/2^k)])$$

$$< (\sum_{k=l}^{\infty} w_k) \exp(nH(p) + T_n - \frac{1}{2}2^l \log(n/2^l) + \delta_1 2^l \log(n/2^l)).$$

δ_1 is a small positive number, at our disposition. The term $\varepsilon_{k,n}$ in (5.19) was dropped, being negative; see (5.18c). The displayed inequality holds by (5.17) in [8].

Compare (5.48) and (5.50): the early zone and midzone are negligible relative to the late and very late high zones, provided $\delta_1 2^l + \varepsilon < 2\delta_0 2^l / \log 2$. It is the minor bit of algebra in (5.47) and (5.49) that seems to determine the critical rate of decay for the w 's in Theorem (1.6).

The endzone goes away, as usual; the early high zone can also be eliminated. We do not know (or need to know) how much posterior mass is in

the middle high zone. Informally, posterior weight shifts so far out that there are only $O(\log n/n)$ observations per parameter. The argument can be completed as in (5.18) in [8]. ■

REFERENCES

- [1] A. D. Barbour, L. Holst and S. Janson, *Poisson Approximation*, Oxford University Press, 1992.
- [2] D. Blackwell, *Large deviations for martingales*, to appear in the *Le Cam Festschrift*, D. Pollard (Ed.), Springer, New York 1995.
- [3] B. de Finetti, *La probabilità, la statistica, nei rapporti con l'induzione, secondo diversi punti di vista*, Centro Internazionale Matematica Estivo Cremonese, Rome 1959. English translation in: B. de Finetti, *Probability, Induction, and Statistics*, Wiley, New York 1972.
- [4] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, Jones and Bartlett, Boston 1993.
- [5] J.-D. Deuschel and D. W. Stroock, *Large Deviations*, Academic Press, Boston 1988.
- [6] P. Diaconis and D. Freedman, *On the consistency of Bayes estimates (with discussion)*, *Ann. Statist.* 14 (1986), pp. 1–67.
- [7] – *On the uniform consistency of Bayes estimates for multinomial probabilities*, *ibidem* 18 (1990), pp. 1317–1327.
- [8] – *Nonparametric binary regression: a Bayesian approach*, *ibidem* 21 (1993), pp. 2108–2137.
- [9] – *Nonparametric Bayesian regression: A review*, Technical Report No. 408, Department of Statistics, University of California, to appear in the *Le Cam Festschrift*, D. Pollard (Ed.), Springer, New York 1995.
- [10] L. E. Dubins and L. J. Savage, *How to Gamble if You Must: Inequalities for Stochastic Processes*, Dover, New York 1985.
- [11] W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. II, 2nd ed., Wiley, New York 1971.
- [12] R. A. Fisher, *On the interpretation of χ^2 from contingency tables, and the calculation of P*, *J. Roy. Statist. Soc.* 85 (1922), pp. 87–94.
- [13] D. Freedman, *Another note on the Borel–Cantelli lemma and the strong law, with the Poisson approximation as a by-product*, *Ann. Probab.* 6 (1973), pp. 910–925.
- [14] L. Holst, *Two conditional limit theorems with applications*, *Ann. Statist.* 7 (1979), pp. 551–557.
- [15] V. Kolchin, B. Sevastyanov and V. Chistyakov, *Random Allocations*, Wiley, New York 1978.
- [16] P. S. de Laplace, *Mémoire sur la probabilité des causes par les événements*, *Mémoires de mathématique et de physique présentés à l'Académie Royale des Sciences, par divers savants, et lus dans ses assemblées* 6 (1774). Reprinted in Laplace's *Oeuvres Complètes* 8, pp. 27–65. English translation by S. Stigler, *Statist. Sci.* 1 (1986), pp. 359–378.
- [17] V. Petrov, *Sums of Independent Random Variables*, Springer, New York 1975.
- [18] D. Siegmund, *Sequential Analysis: Tests and Confidence Intervals*, Springer, New York 1985.

P. Diaconis, Mathematics Department
Harvard University
Cambridge, Massachusetts 02138, U.S.A.

D. Freedman, Statistics Department
University of California
Berkeley, California 94720, U.S.A.

Received on 26.11.1993

