

RIDGE TYPE ESTIMATORS OF
MULTINOMIAL CELL PROBABILITIES¹

Ayodele Ighodaro and Thomas Santner

School of Operations Research and Industrial Engineering
Cornell University
Ithaca, New York, U.S.A.

I. INTRODUCTION AND SUMMARY

Let $\underline{W} = (W_1, \dots, W_t)$ have a multinomial distribution based on N (known) trials and having unknown vector of cell probabilities $\underline{p} \in S \equiv \{q \in R^t: q_i \geq 0 \forall i \text{ and } \sum q_i = 1\}$. When the range of a summation or product is from 1 to t it will be suppressed for ease of notation. This paper develops analogues of normal theory ridge regression estimators for the problem of simultaneously estimating \underline{p} . Sections 1 to 4 will study the case of squared error loss (SEL), $L_S(\underline{p}, \underline{a}) \equiv N \sum (p_i - a_i)^2 = N \|\underline{p} - \underline{a}\|^2$, because of its wide use (see Bishop, Fienberg, and Holland [5], Chapter 12 and the references therein), but other loss functions will be mentioned in Section 5 because of their ability to penalize zero guesses of positive p_i .

The maximum likelihood estimator (mle) $\hat{\underline{p}} \equiv \underline{W}/N$ is known to be admissible under SEL for all t (Johnson [21]; Alam [3]; Brown [6]); hence there is no Stein-effect for this problem. Furthermore, there is a unique admissible minimax estimator (Steinhaus [25]; Trybula [29]). In the absence of any other information about \underline{p} one of these estimators might be satisfactory. However, $\hat{\underline{p}}$ has risk $R_S(\underline{p}, \hat{\underline{p}}) \equiv E_{\underline{p}} [L_S(\underline{p}, \hat{\underline{p}})] = 1 - \|\underline{p}\|^2$ which is near zero only when \underline{p} is near a vertex. This led Good [13], [14],

¹This research was supported by the National Science Foundation under Grant No. ENG-7906914.

Sutherland, Fienberg and Holland [27], Albert [4] and others to consider Bayesian motivated estimators of \underline{p} which dominate $\hat{\underline{p}}$ over a "large" portion of S . In particular, it is well known that the unique Bayes estimator of \underline{p} versus the conjugate Dirichlet prior with mean $E[\underline{p}] = \underline{\lambda}$ and variance-covariance matrix $E[(\underline{p}-\underline{\lambda})(\underline{p}-\underline{\lambda})'] = (D(\underline{\lambda})-\underline{\lambda}\underline{\lambda}')/(K+1)$ is

$$(1.1) \quad \hat{\underline{p}}_K = (W+K\underline{\lambda})/(N+K) = \omega \hat{\underline{p}} + (1-\omega)\underline{\lambda}$$

where $K > 0$ and $\underline{\lambda} = (\lambda_1, \dots, \lambda_t) \in S$ are known, $D(\underline{\lambda})$ is diagonal with the elements of $\underline{\lambda}$ and $\omega = N/(N+K)$. Here and throughout the paper vectors are column vectors and prime denotes transpose. Formally $K = 0$ make sense in (1.1) and yields the mle $\hat{\underline{p}}_0 \equiv \hat{\underline{p}}$ of \underline{p} .

Section 2 shows that the class of estimators $\{\hat{\underline{p}}_K: K \geq 0\}$ is the analogue of the class of ridge estimators in a number of frequentist senses as well as in the well-known Bayesian sense. Section 3 uses the properties of $\hat{\underline{p}}_K$ developed in Section 2 to illustrate the construction of two ridge estimators; it is possible to construct other ridge analogues. Small sample simulation studies are presented in Section 4 which compare the current estimators with those previously proposed; it concentrates on the large sparse multinomial framework introduced by Fienberg and Holland [11]. Section 5 critiques two other loss functions; it summarizes some possible approaches for estimating \underline{p} under various model assumptions in these cases.

II. RIDGE ESTIMATION

It is well known that the class of ridge estimators can be developed from a Bayesian viewpoint by postulating $\underline{Y} = \underline{X}\underline{\beta} + \underline{\epsilon}$ for an $n \times 1$ vector of data \underline{Y} where \underline{X} is a known $n \times p$ matrix of rank p , $\underline{\beta}$ is a $p \times 1$ vector of unknown parameters and $\underline{\epsilon}$ is an $n \times 1$ vector of

experimental errors satisfying $\varepsilon \sim N_n(0, \sigma^2 I_n)$ and by assuming

$\beta \sim N_p(0, \frac{\sigma^2}{K} I_p)$. The Bayes estimator of β with respect to squared error loss is well known to be

$$(2.1) \quad \hat{\beta}_{\hat{\nu}K} = (X'X + KI_p)^{-1} X'Y.$$

Formally $\hat{\beta}_{\hat{\nu}0} = (X'X)^{-1} X'Y$ also makes sense in (2.1); $\hat{\beta}_{\hat{\nu}0}$ is, of course, the mle and BLUE of β . $\hat{\beta}_{\hat{\nu}0}$ is inadmissible under squared error loss when $p \geq 3$; its (summed) mean squared length satisfies

$$(2.2) \quad E_{\beta, \sigma^2} [||\hat{\beta}_{\hat{\nu}0}||^2] - ||\beta||^2 = \sigma^2 \sum_{j=1}^p \lambda_j^{-1}$$

where $\lambda_1, \dots, \lambda_p$ are the eigenvalues of $X'X$. For any $K > 0$, $\hat{\beta}_{\hat{\nu}K}$ is biased; however, its length, $||\hat{\beta}_{\hat{\nu}K}||^2$, is shorter than that of $\hat{\beta}_{\hat{\nu}0}$ (see P1 below).

Hoerl and Kennard [18] developed the following properties and characterizations of $\hat{\beta}_{\hat{\nu}K}$.

P1: The distance $||\hat{\beta}_{\hat{\nu}K}|| = (\hat{\beta}_{\hat{\nu}K}' \hat{\beta}_{\hat{\nu}K})^{1/2}$ of $\hat{\beta}_{\hat{\nu}K}$ to (the prior mean) 0 ($= E[\beta]$) is a continuous monotone decreasing function of K such that $||\hat{\beta}_{\hat{\nu}K}|| \rightarrow 0$ as $K \rightarrow \infty$.

P2: $\hat{\beta}_{\hat{\nu}K}$ is a restricted maximum likelihood (least squares) estimator of β . Denote the residual sum of squares of β by $\psi(\beta) = ||Y - X\beta||^2$; for fixed $K > 0$, $\hat{\beta}_{\hat{\nu}K}$ minimizes $\psi(\beta)$ among β in the sphere $B_K \equiv \{\beta \in R^p: ||\beta|| \leq ||\hat{\beta}_{\hat{\nu}K}||\}$. Hence $\psi(\hat{\beta}_{\hat{\nu}K})$ is increasing in K by P1.

P3: If $||\beta||^2$ is bounded then there exists a $K_0 > 0$ such that

$$E_{\beta} [||\hat{\beta}_{\hat{\nu}K} - \beta||^2] \leq E_{\beta} [||\hat{\beta}_{\hat{\nu}0} - \beta||^2] \text{ for } 0 < K < K_0.$$

P2 characterizes $\hat{\beta}_{\hat{\nu}K}$ while P3 suggests the possibility of constructing minimax ridge estimators by using adaptive (stochastic)

is concave in δ_i and is therefore maximized when

$$(2.5) \quad \frac{1}{K} \frac{W_i}{\delta_i} + \frac{\lambda_i}{\delta_i} = \frac{N+K}{K} \quad (1 \leq i \leq t) \Leftrightarrow \delta_i = \hat{p}_{iK} \quad (1 \leq i \leq t).$$

Substituting (2.5) into the right hand side of (2.4) gives

$$\frac{1}{K} \sum W_i \ln \delta_i \leq \frac{1}{K} \sum W_i \ln \hat{p}_{iK}$$

for all $\delta \in B_K^e$ and completes the proof.

The analogue of a more severe form of P3 is immediate.

THEOREM 2.3. *For any $p \in S$ with at least two non-zero components, the mean square entropy distance of $\hat{p}_{\lambda K}$ to p for any $K > 0$ is less than that of the mle \hat{p} .*

Remark 2.1. Unlike the regression analog P3, there is no upper bound on the value of K which "improves" $\hat{p}_{\lambda K}$ over \hat{p} .

Remark 2.2 If λ has one or more zero components then Theorems 2.1 to 2.3 still hold though their proofs must be modified.

Remark 2.3. Marquardt [22] proves that $\hat{\beta}_{\lambda K}$ is the maximum likelihood estimator of β when the data Y are supplemented by fictitious data $Y^* \equiv 0$ drawn from an experiment orthogonal to the original design. An analogous statement is true for $\hat{p}_{\lambda K}$; $\hat{p}_{\lambda K}$ is the maximum likelihood estimator of p when W is supplemented by $W^* = (K\lambda_1, \dots, K\lambda_t)$.

III. RIDGE ESTIMATORS OF p

The admissibility of \hat{p} insures that unlike the regression case it is impossible to find an estimator \hat{K} of K so that $\hat{p}_{\hat{K}}$ dominates \hat{p} . Yet it is still reasonable to study analogues of ridge estimators in the hope they have reasonable mean squared

error properties relative to \hat{p}_{λ} for p "away" from the vertices. So now assume $K > 0$ but unknown while λ is known. A differentiation shows (see [5], page 407-8) that the K which minimizes $R_S(p, \hat{p}_{\lambda K})$ is

$$(3.1) \quad K^* = K^*(p, \lambda) \equiv (1 - ||p||^2) / ||p - \lambda||^2.$$

We consider two iterative estimators of K^* or equivalently of $\omega^* = N / (N + K^*)$. The idea is to use the current estimate of ω^* to obtain an improved estimate of p which in turn can be used to re-estimate ω^* .

Define the sequence $\{\omega_i\}_{i=1}^{\infty}$ by the algorithm:

1. Make an initial estimate ω_0 of ω^* , set $i=1$ and go to Step 2.
2. Set $\delta_{\lambda i} = \omega_{i-1} \hat{p} + (1 - \omega_{i-1}) \lambda$ and go to Step 3.

$$3. \quad \text{Set } \omega_i = \frac{N ||\delta_{\lambda i} - \lambda||^2}{N ||\delta_{\lambda i} - \lambda||^2 + 1 - ||\delta_{\lambda i}||^2}, \text{ increment } i \text{ by } 1, \text{ and go to}$$

Step 2.

Set $\omega_{ST} = \lim_{i \rightarrow \infty} \omega_i$ and $\hat{p}_{\lambda ST} = \omega_{ST}^* \hat{p} + (1 - \omega_{ST}^*) \lambda$ when the limit exists.

When $W = N\lambda$ then it is easy to check that $\omega_{ST}^* = 0$ and $\hat{p}_{\lambda ST} = \lambda$.

When $W = Ne_i$ where $e_i = (0, \dots, 0, 1, 0, \dots, 0)$ is the i^{th} unit vector then it is also obvious that $\omega_{ST}^* = 1$ and $\hat{p}_{\lambda ST} = e_i$ for $1 \leq i \leq t$. The value of ω_{ST}^* for arbitrary W is given in Theorem 3.1. First an elementary lemma will be stated.

LEMMA 3.1. Suppose $W \neq \lambda$ and α and v are defined as

$\alpha \equiv \lambda' (\hat{p} - \lambda) / N ||\hat{p} - \lambda||^2$ and $v \equiv (1 - \lambda' \lambda) / N ||\hat{p} - \lambda||^2$, respectively then $v \geq 2\alpha$.

Proof. $v \geq 2\alpha \Leftrightarrow 1 + \lambda' \lambda - 2\lambda' \hat{p} \geq 0$
 $\Leftrightarrow (\lambda - \hat{p})' (\lambda - \hat{p}) + (1 - \hat{p}' \hat{p}) \geq 0$

which is obvious.

THEOREM 3.1. Define $\gamma = (N-1)/N$ and α, ν as in Lemma 3.1. When the iterative estimator of the optimal constant ω^* begins at Step 1 with $\omega_0 = 1$ then the sequence $\{\omega_i\}_{i=1}^{\infty}$ converges to the limit

$$\omega_{ST}^* = \begin{cases} 0 & \text{when either } (2\alpha+1)^2 < 4\gamma\nu \text{ or} \\ & (2\alpha+1) < 0 \text{ or} \\ & (2\gamma < (2\alpha+1) \text{ and } p \notin \{e_1, \dots, e_t\}) \\ \frac{(2\alpha+1) + [(2\alpha+1)^2 - 4\gamma\nu]^{1/2}}{2\gamma} & \text{when } ((2\alpha+1)^2 \geq 4\gamma\nu \\ & \text{and } 2\gamma \geq (2\alpha+1) \geq 0) \\ 1 & \text{otherwise.} \end{cases}$$

Proof. After substitution and some algebra it can be shown that

$$\omega_i = \frac{N \left| \frac{\delta_i - \lambda}{\lambda} \right|^2}{N \left| \frac{\delta_i - \lambda}{\lambda} \right|^2 + 1 - \left| \frac{\delta_i}{\lambda} \right|^2} = f(\omega_{i-1}),$$

say, where $f(\omega) = \omega^2 / (\gamma\omega^2 - 2\alpha\omega + \nu)$. Here $f: [0,1] \rightarrow [0,1]$; differentiation gives $f'(\omega) = 2\omega(\nu - \alpha\omega) / (\gamma\omega^2 - 2\alpha\omega + \nu)^2 > 0$ for $1 \geq \omega > 0$ by Lemma 3.1. Thus $\{\omega_i\}_{i=1}^{\infty}$ is a nonincreasing sequence bounded below and therefore must have a limit. Furthermore $\omega_{i+1} < \omega_i$ whenever $\omega_i > 0$. The limit points of $\{\omega_i\}_{i=1}^{\infty}$ must solve the fixed point equation

$$(3.2) \quad \begin{aligned} f(\omega) &= \omega \\ \Leftrightarrow \omega(\gamma\omega^2 - (2\alpha+1)\omega + \nu) &= 0 \\ \Leftrightarrow \omega \in \{\omega^1, \omega^2, \omega^3\}, \end{aligned}$$

$$\omega^1 \equiv 0, \quad \omega^2 \equiv [(2\alpha+1) + [(2\alpha+1)^2 - 4\gamma\nu]^{1/2}] / 2\gamma, \quad \text{and}$$

$$\omega^3 \equiv [(2\alpha+1) - [(2\alpha+1)^2 - 4\gamma\nu]^{1/2}] / 2\gamma.$$

If $(2\alpha+1)^2 < 4\gamma\nu$ then ω^1 is the only fixed point of $f(\omega)$ and

hence $\omega_{ST}^* = 0$. Suppose $(2\alpha+1)^2 \geq 4\gamma\nu$; it is easy to check that $g(\omega) = \gamma\omega^2 - (2\alpha+1)\omega + \nu$ has either both roots positive or both roots negative since $g(0) = \nu > 0$. Straightforward algebra shows that $\omega^2 \geq 0, \omega^3 \geq 0 \Leftrightarrow (2\alpha+1) \geq 0$. Hence if $(2\alpha+1) < 0$ then $\omega^2 \leq 0$ and $\omega^3 \leq 0 \Rightarrow \omega^1 = 0$ is again the only fixed point of $f(\omega)$ to which $\{\omega_i\}_{i=1}^\infty$ can converge since for every i , $\omega_i \in [0,1] \Rightarrow \omega_{i+1} = f(\omega_i) \geq f(0) = 0$. If $(2\alpha+1) \geq 0$ and $2\gamma - (2\alpha+1) \geq 0$ then straightforward but tedious algebra shows

$$\omega^2 \leq 1 \Leftrightarrow \hat{p}_\lambda, \hat{p}'_\lambda \leq 1$$

which is always true. Hence $1 \geq \omega^2 \geq \omega^3 \geq \omega^1 = 0$ and since $\omega_0 = 1$ then for every i , $\omega_i \in [\omega^2, 1] \Rightarrow \omega_{i+1} = f(\omega_i) \geq f(\omega^2) = \omega^2$ and thus $\omega_{ST}^* = \omega^2$. The remaining subcases proceed along similar lines and the proof is completed.

Remark 3.1. $2\alpha+1 \geq 0 \Leftrightarrow \frac{\lambda'(\hat{p}_\lambda - \lambda)}{N \|\hat{p}_\lambda - \lambda\|} \geq -\frac{1}{2}$

$$\Leftrightarrow (\hat{p}_\lambda - \lambda)' \left(\hat{p}_\lambda - \frac{N-2}{N} \lambda \right) \geq 0$$

and hence for most λ is will be the case that $(2\alpha+1) \geq 0$.

Remark 3.2. The proof shows that while ω_{ST}^* always exists, its value can depend on the starting value ω_0 . For example, if $(2\alpha+1) \geq 0$ and $2\gamma \geq (2\alpha+1) \geq 0$ then $1 \geq \omega^2 \geq \omega^3 \geq \omega^1 = 0$ hence if $\omega_0 \in [\omega^3, \omega^2)$ then $\omega_{ST}^* = \omega^3$ rather than ω^2 . A similar phenomenon was noted in the ridge regression case by Hocking, Speed and Lynn [17].

It is instructive to consider the special case $\lambda = \xi \equiv (1/t, \dots, 1/t)$, it is easy to see $\alpha = 0$ and $\nu = (t-1)/\chi^2$ where $\chi^2 = \frac{t}{N} \sum_{i=1}^t (X_i - N/t)^2$. When $N > 2$ Theorem 3.1 simplifies to

$$\omega_{ST}^* = \begin{cases} 0, & \chi^2 < 4 \frac{(N-1)}{N} (t-1) \\ \frac{N[1 + \{1 - 4(\frac{N-1}{N}) \frac{t-1}{\chi^2}\}^{1/2}]}{2(N-1)}, & \text{otherwise} \end{cases}$$

since $(2\alpha+1) = 1 > 0$ and $2\gamma - (2\alpha+1) = 2\gamma-1 = (N-2)/N > 0$; hence $\hat{p}_{\lambda ST} = 0 \cdot \hat{p} + (1-0) \cdot \hat{c} = \hat{c}$ when $\chi^2 \leq 4(t-1)(N-1)/N$. We conclude that it requires an overwhelming amount of evidence against cell homogeneity to estimate the vector \hat{p} by anything other than \hat{c} . Furthermore when $\chi^2 \geq 4(t-1)(N-1)/N$ then $\omega_{ST}^* > N/2(N-1) > 1/2$ thus \hat{p}_{ST} has a large discontinuity. Again a similar problem is encountered in the normal means case; see Hemmerle [16] and Hocking, Speed and Lynn [17].

The preceding discussion suggests that ω_{ST}^* may iterate too far; one possible remedy is to use a finite step estimator of ω^* . If $\omega_0 = 1$ then the one step estimator of ω^* is ω_1 which is the maximum likelihood estimator of ω^* ; denote the corresponding estimator of \hat{p} by $\hat{p}_{\lambda ML}$. This estimator has been studied extensively beginning with Fienberg and Holland [10].

This paper will compare $\hat{p}_{\lambda ML}$ with the estimator which iterates one additional step; denote the two step estimator of ω^* by $\omega_{TS}^* \equiv \omega_2$ (initialized at $\omega_0 = 1$). After some algebra it can be shown that

$$\begin{aligned} \omega_{TS}^* &= f(\omega_2) = f(f(1)) \\ &= [\gamma - 2\alpha(\gamma - 2\alpha + \nu) + \nu(\gamma - 2\alpha + \nu)^2]^{-1} \end{aligned}$$

where α , ν and γ are defined above. When $\lambda = \hat{c}$ then ω_{TS}^* simplifies to

$$\begin{aligned} \omega_{TS}^* &= \left[\frac{N-1}{N} + \frac{t-1}{\chi^2} \left(\frac{N-1}{N} + \frac{t-1}{\chi^2} \right)^2 \right]^{-1} \text{ and} \\ \hat{p}_{\lambda TS} &= \omega_{TS}^* (\hat{p}_{\lambda} - \hat{c}) + \hat{c}. \end{aligned}$$

Remark 3.3. In contrast, when $\lambda = \xi$ then the one step estimator is

$$\hat{p}_{\lambda ML} = \left[\frac{N-1}{N} + \frac{t-1}{\chi^2} \right]^{-1} (\hat{p}_{\lambda} - c) + \xi.$$

When N is large $\hat{p}_{\lambda ML}$ is approximately the same as

$$\hat{p}_{\lambda G} = \left[1 + \frac{t-1}{\chi^2} \right]^{-1} (\hat{p}_{\lambda} - c) + \xi.$$

Good [13] found $\hat{p}_{\lambda G}$ to be a reasonable approximation in several numerical examples to the so called "type II" maximum likelihood (or empirical Bayes) estimator of p_{λ} .

An analogue of the McDonald and Galarneau [23] method for estimating K will now be developed. Let $\widehat{||p||}_e^2$ be an estimator of $||p||_e^2$; choose K so that

$$(3.3) \quad ||\hat{p}_{\lambda K}||_e^2 = \widehat{||p||}_e^2.$$

Unlike the regression case, an unbiased estimate of $||p||_e^2$ does not exist; we make use of a concept due to Haldane [15].

Definition 3.1. The estimator \hat{q} of the real function $q = q(\theta)$ of the parameter θ is said to be *almost unbiased* if

$$E_{\theta} [\hat{q}] = q(\theta) + o(n^{-\lambda})$$

for all θ and for all $\lambda \in \{1, 2, \dots\}$.

Using a result of Cook, Kerridge, and Pryce [8] an almost unbiased estimator of $||p||_e^2$ is

$$(3.4) \quad \widehat{||p||}_e^2 = D(N) + \sum \lambda_i (\lambda n \lambda_i^{-D(W_i)})$$

where $D(y) = \frac{d \ln \Gamma(y+1)}{dy}$ is the digamma function. Define the fitted length estimator to be

$$\hat{p}_{\lambda_{FL}} = \omega_{FL}^* \hat{p} + (1 - \omega_{FL}^*) \lambda$$

where $\omega_{FL}^* = N / (N + K_{FL}^*)$ and

$$K_{FL}^* = \begin{cases} 0, & \widehat{\|p\|}_e^2 \geq \|\hat{p}\|_e^2 \\ \infty, & \widehat{\|p\|}_e^2 \leq 0 \\ \text{solution of (3.3), otherwise.} \end{cases}$$

Remark 3.4. $\hat{p}_{\lambda_{FL}} > 0$ for all \mathbb{W} since if $W_j = 0$ for some j then $\|\hat{p}\|_e^2 = \infty > \widehat{\|p\|}_e^2$ from (3.4).

Remark 3.5. When $\lambda = \zeta$ then K_{FL}^* simplifies to

$$K_{FL}^* = \begin{cases} 0, & [\frac{1}{t} \sum D(W_i)] - D(N) \leq \frac{1}{t} \sum \ln(W_i/N) \\ \infty, & [\frac{1}{t} \sum D(W_i)] - D(N) \geq \frac{1}{t} \sum \ln(1/t) \\ \text{solution of (3.5), otherwise} \end{cases}$$

where

$$(3.5) \quad [\frac{1}{t} \sum D(W_i)] - D(N) = \frac{1}{t} \sum \ln(\hat{p}_{ik})$$

$$\text{and } \hat{p}_{\lambda_K} = (\hat{p}_{1K}, \dots, \hat{p}_{tK}).$$

The next section compares the risk properties of the estimators proposed above with other estimators from the literature.

IV. SMALL SAMPLE SITUATIONS

First, some estimators proposed by Sutherland et al. [27] will be reviewed, their theoretical large sparse multinomial asymptotic behavior will be stated and finally the results of a simulation experiment to compare their small sample risks with those of Section 3 will be presented.

4.1 Review of Previous Estimators

Sutherland et al. proposed two methods of estimating the optimal constant K^* in equation (3.1). The first method was mentioned in Section 3; it estimates K^* by $K_{ML}^* = (N^2 - \sum W_i^2) / (\sum W_i - N\lambda_i)^2$, the mle of K^* based on the conditional distribution of \mathbb{W} given \mathbb{p} . The second ("ratio unbiased") method uses unbiased estimators of the numerator and denominator of K^* ,

$$K_{RU}^* = (N^2 - \sum W_i^2) / (\sum W_i^2 - 2(N-1) \sum W_i \lambda_i + N(N-1) \sum \lambda_i^2 - N).$$

Denote the estimators of \mathbb{p} corresponding to K_{ML}^* and K_{RU}^* by $\hat{\mathbb{p}}_{ML}$ and $\hat{\mathbb{p}}_{RU}$, respectively.

Sutherland [26] analyzes the asymptotic behavior of the risk of a large class of estimators, $C(\lambda)$, as $N \rightarrow \infty$, $t \rightarrow \infty$ and $\rho \equiv N/t$ fixed (large, sparse tables); $C(\lambda)$ includes $\hat{\mathbb{p}}_{ML}$, $\hat{\mathbb{p}}_{RU}$ and $\hat{\mathbb{p}}$. Since the dimension of the true probability vector \mathbb{p} and of λ increases to infinity he fixes two densities $p(u)$ and $\lambda(u)$ on $[0,1]$ and defines

$$p_i = \frac{1}{t} p\left(\frac{i-.5}{t}\right), \quad i = 1, \dots, t \text{ and}$$

$$\lambda_i = \frac{1}{t} \lambda\left(\frac{i-.5}{t}\right), \quad i = 1, \dots, t.$$

for integer t . Sutherland proves that to $o(1)$ the leading term of $R(\mathbb{p}, \hat{\mathbb{p}}_{RU})$ is at least as small as the leading term of any $\mathbb{p}^* \in C(\lambda)$; in particular it is smaller than the leading term of $R(\mathbb{p}, \hat{\mathbb{p}}_{ML})$ which is smaller than the leading term of $R(\mathbb{p}, \hat{\mathbb{p}})$. In addition he computes the asymptotic risk to $o(t^{-1})$ for arbitrary $\mathbb{p}^* \in C(\lambda)$. See also Bishop, Fienberg and Holland [5], Chapter 12.

4.2 Small Sample Studies

Our simulations attempt to answer three questions:

- (1) How large must N and t be before the large sparse asymptotics are correct?

- (2) Does the large sample improvement of $\hat{p}_{\lambda_{RU}}$ over $\hat{p}_{\lambda_{ML}}$ and/or $\hat{p}_{\lambda_{ML}}$ over \hat{p}_{λ} carry over to small samples?
- (3) How do the risk characteristics of the estimators of Sutherland et al. [27] compare with those of Section 3?

Remark 4.1. Bishop, Fienberg, and Holland [5] suggest improving $\hat{p}_{\lambda_{RU}}$ by setting K_{RU}^* equal to zero whenever it would be negative; we denote the corresponding estimator of p_{λ} by $\hat{p}_{\lambda_{RUB}}$. A modification that appears more reasonable is to set K_{RU}^* equal to $+\infty$ whenever it would be negative. To see this, note that K_{RU}^* is negative iff the unbiased estimate of $||p-\lambda||^2$ is negative. Thus setting the estimate of $||p-\lambda||^2$ equal to zero will result in our modification. Denote the corresponding estimator of p_{λ} by $\hat{p}_{\lambda_{RUI}}$. Thus the simulation study examines the following estimators:

- (1) \hat{p}_{λ} , (2) $\hat{p}_{\lambda_{RUB}}$, (3) $\hat{p}_{\lambda_{RUI}}$, (4) $\hat{p}_{\lambda_{ML}}$, (5) $\hat{p}_{\lambda_{ST}}$, (6) $\hat{p}_{\lambda_{TS}}$, and (7) $\hat{p}_{\lambda_{FL}}$.

Throughout this study $\lambda(u)$ was fixed to be a uniform distribution over $[0,1]$, i.e. $\lambda = c = (t^{-1}, \dots, t^{-1})$ for each t . In this case all seven estimators have invariant risk with respect to permutations of p_{λ} and hence it suffices to study densities $p(u)$ which are monotone increasing. The beta densities

$$p(u|\eta) \equiv \begin{cases} \eta u^{\eta-1}, & 0 \leq u \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

for $\eta = 1, 1.4, 3.5$ and 10 were selected to generate the p_{λ} sequences; for each (t, η) , p_i was set equal to $p_i^*/\sum p_i^*$ where $p_i^* \equiv t^{-1}p((i-.5)t^{-1}|\eta)$ for $i = 1, \dots, t$. The vector p_{λ} moves from the center of the simplex at $\eta = 1$ to a point near the vertex at $\eta = 10$. The simulation studied $t = 3, 5(5)20, 50, 100$ cells and $\rho = N/t = 1, 3, 5, 10$; one hundred replications were performed for most (η, t, ρ) combinations (not all combinations were used for $t = 50$ and 100 because of prohibitive computer expense). Among the checks made on the simulation were comparisons of the

estimated risks of $\hat{p}_{\hat{\nu}}$ with their known values.

Extensive tables of the estimated risk values with their standard errors can be found in Ighodaro [19]. Some of the estimated risk values have been plotted in various ways together with the theoretical risk curves (to $o(t^{-1})$) for $\hat{p}_{\hat{\nu}_{RU}}$ and $\hat{p}_{\hat{\nu}_{ML}}$ in Figures 1 through 6. The theoretical risks were computed from the general expressions in Sutherland [26]. The standard errors of the estimated risk values ranged approximately from .003 to .08 although the majority were in the range .03 to .04.

1. When $(\eta, \rho) = (1, 5)$ then $p_{\hat{\nu}} = \zeta$ and

$$(4.1) \quad R(\zeta, \hat{p}_{\hat{\nu}_{RU}}) = 2/t + o(t^{-1}) \text{ and}$$

$$(4.2) \quad R(\zeta, \hat{p}_{\hat{\nu}_{ML}}) = 2.5 + 0.175/t + o(t^{-1}).$$

2. When $(\eta, \rho) = (1.4, 5)$ then

$$(4.3) \quad R(\hat{p}_{\hat{\nu}}, \hat{p}_{\hat{\nu}_{RU}}) = .3076 + 1.4435/t + o(t^{-1}) \text{ and}$$

$$(4.4) \quad R(\hat{p}_{\hat{\nu}}, \hat{p}_{\hat{\nu}_{ML}}) = .4236 + .1408/t + o(t^{-1}).$$

3. When $(\eta, \rho) = (10, 3)$ then

$$(4.5) \quad R(\hat{p}_{\hat{\nu}}, \hat{p}_{\hat{\nu}_{RU}}) = .9275 - 3.9517/t + o(t^{-1}) \text{ and}$$

$$(4.6) \quad R(\hat{p}_{\hat{\nu}}, \hat{p}_{\hat{\nu}_{ML}}) = .9278 - 4.0349/t + o(t^{-1}).$$

Figures 1, 3, and 5 plot equations (4.1) to (4.6) to $o(t^{-1})$ together with the simulated risk curves for $\hat{p}_{\hat{\nu}_{RUI}}$ and $\hat{p}_{\hat{\nu}_{ML}}$. Figure 1 also includes the simulated risk of $\hat{p}_{\hat{\nu}_{RUB}}$; it is clear from this plot that $\hat{p}_{\hat{\nu}_{RUI}}$ is the correct modification to $\hat{p}_{\hat{\nu}_{RU}}$. Hence $\hat{p}_{\hat{\nu}_{RUB}}$ is not considered in Figures 2-6. It can also be concluded that the asymptotic risk expressions to $o(t^{-1})$ are approximately correct for $t \geq 30$.

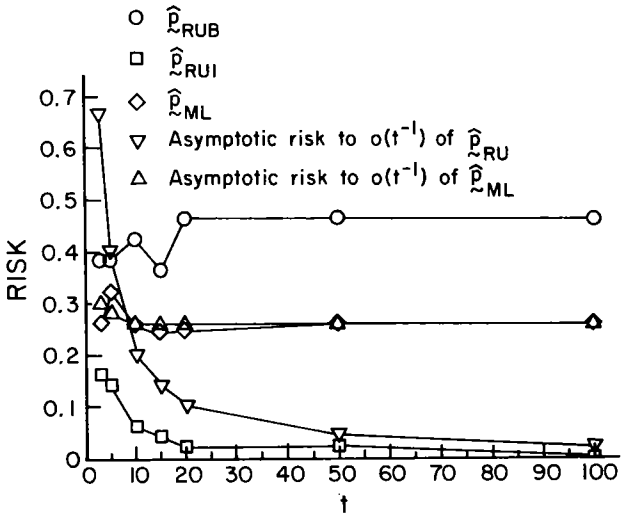


Figure 1. Simulated and asymptotic risks for $(n, \rho) = (1, 5)$.

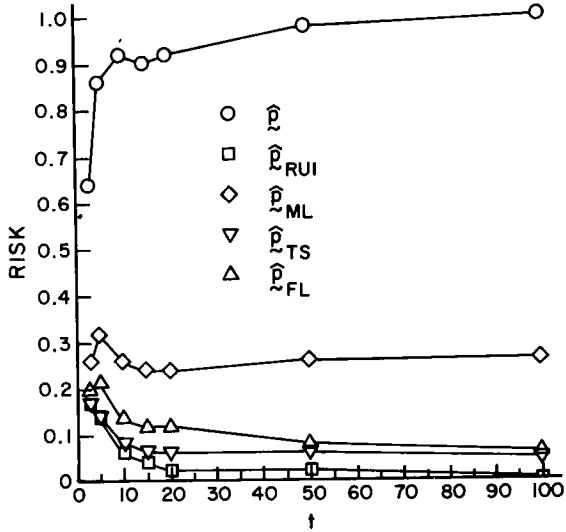


Figure 2. Simulated risk values for $(n, \rho) = (1, 5)$.

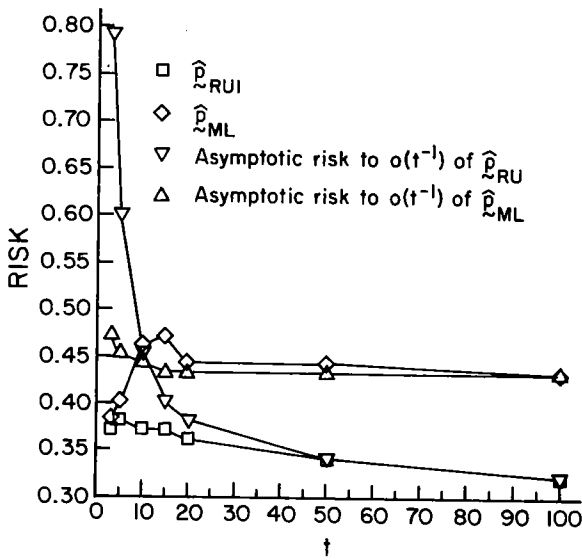


Figure 3. Simulated and asymptotic risks for $(n, \rho) = (1.4, 5)$.

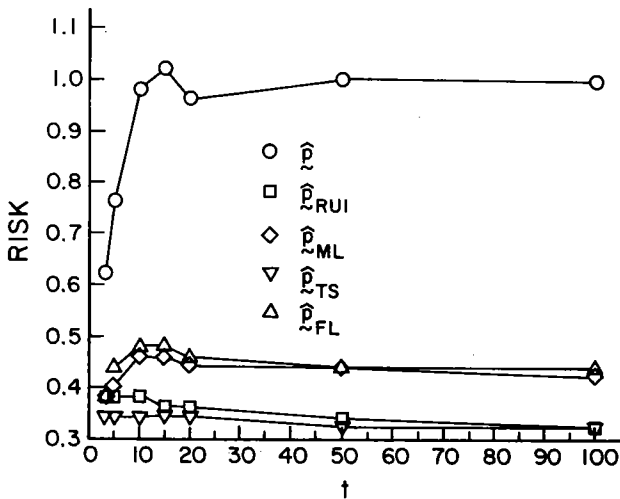


Figure 4. Simulated risk values for $(n, \rho) = (1.4, 5)$.

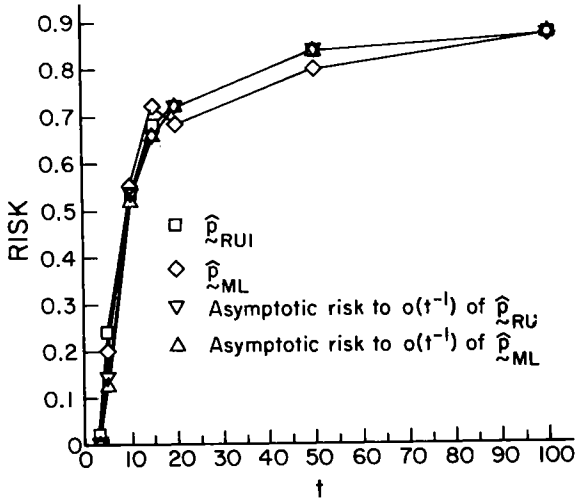


Figure 5. Simulated and asymptotic risks for $(n, \rho) = 10, 3$.

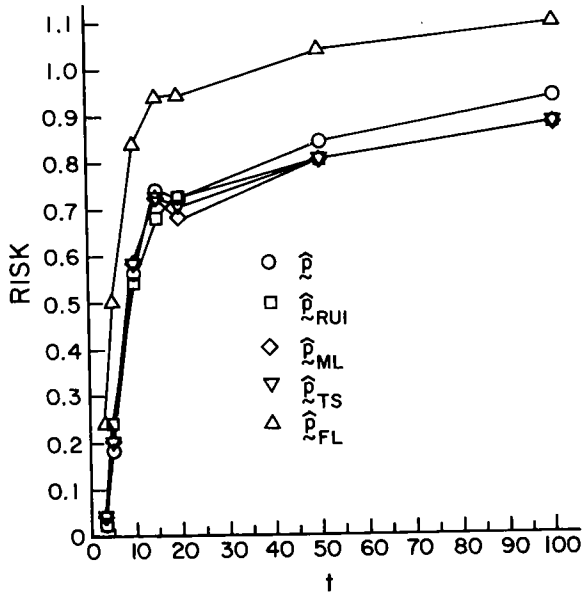


Figure 6. Simulated risk values for $(n, \rho) = (10, 3)$.

Figures 2, 4 and 6 are plots of the simulated risk curves for $\hat{p}_{\lambda_{RUI}}$, $\hat{p}_{\lambda_{ML}}$, \hat{p}_{λ} , $\hat{p}_{\lambda_{TS}}$ and $\hat{p}_{\lambda_{FL}}$. The curve for $\hat{p}_{\lambda_{ST}}$ is not plotted because, as expected, it performs badly over a wide portion of the parameter space. For the remaining estimators we conclude that:

- (a) $\hat{p}_{\lambda_{RUI}}$ and $\hat{p}_{\lambda_{TS}}$ are comparable and better overall than any of the other estimators,
- (b) $\hat{p}_{\lambda_{FL}}$ is superior to the mle \hat{p}_{λ} except for sparse data ($\rho = 3$) near vertices ($\eta = 10$),
- (c) all other estimators significantly outperform \hat{p}_{λ} except near the vertices ($\eta = 10$) where they are comparable to it.

V. DISCUSSION

Despite the favorable risk comparisons of most of the estimators in Section 4 relative to $\hat{p}_{\lambda} = W/N$, it is unknown whether any of them are admissible for N and t fixed. Unique Bayes (hence admissible) estimators with data dependent K do arise from compound Dirichlet priors (Good [13], [14]), but they appear to be analytically intractable. More directly Brown's [6] characterization of the admissible rules might be useful for this problem.

SEL has been adopted throughout this paper because of its widespread use. However, under L_S the locus of equivalent actions \hat{p}_{λ} versus a given $p_i > 0$ is spherical; hence L_S does not differentiate between positive and zero guesses of $p_i > 0$. One alternative, relative squared error loss (RSEL), $L_R(p_{\lambda}, a) \equiv \sum (p_i - a_i)^2 p_i$, deals with a related question. RSEL does not allow positive guesses of zero p_i however it still does allow zero guesses of positive p_i . Olkin and Sobel [24] apply the divergence theorem to show that the mle \hat{p}_{λ} is admissible unique minimax. Alternatively this can be proved by noting the unique Bayes estimator of p_{λ} relative to the Dirichlet prior

$$h(p_{\lambda}) = \Gamma(K) \prod \{p_i^{K\lambda_i - 1} / \Gamma(K\lambda_i)\}, \quad p_{\lambda} \in S$$

is for $1 \leq i \leq t$:

$$\hat{p}_i^B = \begin{cases} 0, & W_i = 0 \\ (W_i + K\lambda_i - 1)/(N+k-t), & W_i \geq 1; \end{cases}$$

in particular $\hat{p}_{\tilde{\nu}}$ is unique Bayes versus the prior with $\tilde{\lambda} = (1/t, \dots, 1/t)$ and $K = t$. Minimacity follows from the fact that $\hat{p}_{\tilde{\nu}}$ has constant risk. Ighodaro [19] develops an asymptotic ($N \rightarrow \infty$, t fixed) application of the James-Stein estimator to the problem of estimating $p_{\tilde{\nu}}$ under the model assumption $p_{\tilde{\nu}} \in T \subset S$ where T is a smooth surface (e.g. $\ln p_{\tilde{\nu}}$ satisfies a log linear model). The small sample risk performance of these estimators are unknown.

Entropy loss (EL), $L_E(p_{\tilde{\nu}}, a_{\tilde{\nu}}) \equiv N \sum p_i \ln(p_i/a_i)$ where $0 \ln 0 = 0$, is an alternative to SEL which effectively differentiates between positive and zero guesses of $p_i > 0$ (see Akaike [1], [2]); if $a_i = 0$ for any $p_i > 0$ then $L_E(p_{\tilde{\nu}}, a_{\tilde{\nu}}) = +\infty$. The mle has entropy risk $R_E(p_{\tilde{\nu}}, \hat{p}_{\tilde{\nu}}) = 0$ or $+\infty$ according as $p_{\tilde{\nu}}$ is a vertex or not. The mle is more dramatically unsatisfactory under L_E than under L_S ; $\hat{p}_{\tilde{\nu}}$ is admissible when the problem has parameter space S but it is inadmissible when the problem has parameter space $S' \equiv \{p_{\tilde{\nu}} \in S: p_{\tilde{\nu}} \text{ not a vertex}\}$ (Ighodaro, Santner and Brown [20]). A theory of empirical or pseudo Bayes rules could be developed for L_E along the lines of Bishop et al. [5].

ACKNOWLEDGMENT

The authors would like to thank Professor M. Todd for suggesting a simplified proof of Theorem 2.1.

REFERENCES

- [1] Akaike, H. (1971). Information theory and an extension of the maximum likelihood principle. *Second International Symposium on Information Theory*, B. Petrov and F. Csaki (eds.). Akademiai, Budapest.

- [2] Akaike, H. (1978). A new look at the Bayes procedure. *Biometrika* 65, 53-59.
- [3] Alam, K. (1979). Estimation of multinomial probabilities. *Ann. Statist.* 7, 282-283.
- [4] Albert, J. (1979). Robust Bayes estimation. Technical Report 79-9, Department of Statistics, Purdue University, W. Lafayette.
- [5] Bishop, Y., Fienberg, S., and Holland, P. (1975). *Discrete Multivariate Analysis*. MIT Press, Cambridge.
- [6] Brown, L. D. (1981). A complete class theorem for statistical problems with finite sample spaces. *Ann. Statist.* 9, 1289-1300.
- [7] Casella, G. (1980). Minimax ridge regression estimation. *Ann. Statist.* 8, 1036-1056.
- [8] Cook, G., Kerridge, D., and Pryce, J. (1974). Estimation of functions of a binomial parameter. *Sankhyā Ser. A* 36, 443-448.
- [9] Draper, N. and Van Nostrand, C. (1979). Ridge regression and James-Stein estimation: review and comments. *Technometrics* 21, 451-466.
- [10] Fienberg, S., and Holland, P. (1970). Methods for eliminating zero counts in contingency tables. *Random Counts on Models and Structures*, G. P. Patel (ed.). Pennsylvania State Univ. Press, University Park.
- [11] Fienberg, S., and Holland, P. (1973). Simultaneous estimation of multinomial cell probabilities. *J. Amer. Statist. Assoc.* 68, 683-691.
- [12] Goel, P. K. and Casella, G. (1976). A note on an explicit solution for generalized ridge regression. Technical Report No. 448, Department of Statistics, Purdue University, West Lafayette.
- [13] Good, I. J. (1965). *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. MIT Press, Cambridge.
- [14] Good, I. J. (1967). A Bayesian significance test for multinomial distributions (with discussion). *J. Roy. Statist. Soc. Ser. B* 29, 399-431.

- [15] Haldane, J.B.S. (1957). Almost unbiased estimates of functions of sequences. *Sankhyā* 17, 201-8.
- [16] Hemmerle, W. J. (1975). An explicit solution for generalized ridge regression. *Technometrics* 17, 309-314.
- [17] Hocking, R., Speed, F., and Lynn, M. (1976). A class of biased estimators in linear regression. *Technometrics* 18, 425-437.
- [18] Hoerl, A., and Kennard, R. (1970). Ridge regression: biased estimation for nonorthogonal problems. *Technometrics* 12, 55-67.
- [19] Ighodaro, A. O. (1980). Ridge and James-Stein methods for contingency tables. Ph.D. Thesis, School of OR/IE, Cornell University, Ithaca.
- [20] Ighodaro, A., Santner, T., and Brown, L. (1980). Some admissibility and complete class results for the multinomial problem under entropy and squared error loss. To appear in *J. Mult. Anal.*
- [21] Johnson, B. M. (1981). On the admissible estimators for certain fixed sample binomial problems. *Ann. Math. Statist.* 42, 1579-1587.
- [22] Marquardt, D. W. (1970). Generalized inverses, ridge regression, biased linear estimation, and nonlinear estimation. *Technometrics* 12, 591-612.
- [23] McDonald, G., and Galarneau, D. (1975). A Monte-Carlo evaluation of some ridge-type estimators. *J. Amer. Statist.* 70, 407-416.
- [24] Olkin, I. and Sobel, M. (1979). Admissible and minimax estimation for the multinomial distribution and for K independent binomial distributions. *Ann. Statist.* 7, 284-290.
- [25] Steinhaus, H. (1957). The problem of estimation. *Ann. Math. Statist.* 28, 633-648.
- [26] Sutherland, M. (1974). Estimation in large sparse multinomials. Ph.D. Thesis, Department of Statistics, Harvard University, Cambridge.
- [27] Sutherland, M., Holland, P., and Fienberg, S. (1974). Combining Bayes and frequency approaches to estimate a multinomial parameter. *Studies in Bayesian Econometrics and Statistics*, S. Fienberg and A. Zellner (eds.). North Holland, Amsterdam.

- [28] Thisted, R. (1976). Ridge regression, minimax estimation, and empirical Bayes methods. Ph.D. Thesis, Technical Report No. 28, Division of Biostatistics, Stanford University, Stanford.
- [29] Trybula, S. (1958). Some problems of simultaneous minimax estimation. *Ann. Math. Statist.* 39, 245-253.