

Achievability of Asymptotic Minimax Regret in Online and Batch Prediction

Kazuho Watanabe

*Graduate School of Information Science
Nara Institute of Science and Technology
8916-5, Takayama-cho, Ikoma, Nara, 630-0192, Japan*

WKAZUHO@IS.NAIST.JP

Teemu Roos

*Department of Computer Science
University of Helsinki
Helsinki Institute for Information Technology HIIT
PO Box 68, FI-00014, Finland*

TEEMU.ROOS@CS.HELSENKI.FI

PETRI.MYLLYMAKI@CS.HELSENKI.FI

Editor: Cheng Soon Ong and Tu Bao Ho

Abstract

The normalized maximum likelihood model achieves the minimax coding (log-loss) regret for data of fixed sample size n . However, it is a batch strategy, i.e., it requires that n be known in advance. Furthermore, it is computationally infeasible for most statistical models, and several computationally feasible alternative strategies have been devised. We characterize the achievability of asymptotic minimaxity by batch strategies (i.e., strategies that depend on n) as well as online strategies (i.e., strategies independent of n). On one hand, we conjecture that for a large class of models, no online strategy can be asymptotically minimax. We prove that this holds under a slightly stronger definition of asymptotic minimaxity. Our numerical experiments support the conjecture about non-achievability by so called last-step minimax algorithms, which are independent of n . On the other hand, we show that in the multinomial model, a Bayes mixture defined by the conjugate Dirichlet prior with a simple dependency on n achieves asymptotic minimaxity for all sequences, thus providing a simpler asymptotic minimax strategy compared to earlier work by Xie and Barron. The numerical results also demonstrate superior finite-sample behavior by a number of novel batch and online algorithms.

Keywords: on-line learning, prediction of individual sequences, asymptotic minimax regret, Bayes mixture, last-step minimax algorithm

1. Introduction

The normalized maximum likelihood (NML) distribution is derived as the optimal solution to the minimax problem which minimizes the worst-case regret in code-length (log-loss) of data with fixed sample size n . Although a direct evaluation of the NML distribution involves the computation of a sum over all possible data sets, taking exponential time, linear-time algorithms have been developed for certain models such as multinomials (Kontkanen and Myllymäki, 2007; Silander et al., 2010). However, the computation of the NML distribution is still intractable for most models.

Approximating the minimax solution by other easily implementable strategies has been studied. Asymptotic minimaxity is a key feature of such strategies, where the worst-case code-length converges to that of the NML as the sample size tends to infinity. For the multinomial model, [Xie and Barron \(2000\)](#) showed that a Bayes procedure defined by a modified Jeffreys prior, where additional mass is assigned to the boundaries of the parameter space, can achieve asymptotic minimax optimality. An alternative technique to this procedure was studied for a more general model class ([Takeuchi and Barron, 1997](#)).

In the context of online prediction of individual sequences, the focus has been on prediction strategies which can be computed without knowing the sequence length n in advance. We call such strategies *online*, while strategies that take advantage of the knowledge of the sample size n are called *batch*. Online strategies are essential in processing large data sets and especially streaming data where the knowledge of n is not available. Theoretical properties of online strategies have intensively been studied in theories of online learning. For online strategies, regret bounds of the form $k \ln n + O(1)$, where k is a constant, have been obtained ([Azoury and Warmuth, 2001](#); [Cesa-Bianchi and Lugosi, 2001](#); [Freund, 1996](#)). Furthermore, it was proved for the Bernoulli model and the exponential families with a constrained parameter space that the minimax optimal regret is achieved, up to the $O(1)$ term, by the Bayesian strategy using the Jeffreys prior and the last-step minimax strategy (a.k.a. the sequential normalized maximum likelihood) ([Takimoto and Warmuth, 2000](#); [Kotłowski and Grünwald, 2011](#)). That is, if the regret of the NML is asymptotically expanded as $k^* \ln n + c^* + o(1)$ with constants k^* and c^* , $k = k^*$ holds for these strategies. The asymptotic minimax optimality examines if the optimal constant c^* is also achieved and the maximum regret matches that of the NML up to the $o(1)$ term.

In this paper, we investigate achievability of asymptotic minimaxity by batch and online strategies. We consider a slightly stronger asymptotic minimax property and prove that under a generic condition on the model class, it cannot be achieved by any online strategy (Thm. 1). We conjecture that a similar result also holds for the standard asymptotic minimax notion. We also show that for the multinomial model, a sample-size-dependent Bayes procedure defined by a simpler prior than the modified Jeffreys prior in [Xie and Barron \(2000\)](#) achieves asymptotic minimaxity under the standard definition, as well as approximately in our stronger sense (Thm. 4). Through numerical experiments (Sect. 4), we demonstrate the achievability of asymptotic minimaxity for batch strategies. We also investigate the behavior of a generalization of the last-step minimax algorithm, which we call the k -last-step minimax algorithm and which is online. We demonstrate that while for large k , its performance is very near asymptotic minimax, it fails to achieve it exactly, in line with our conjecture. Lastly, the numerical results demonstrate superior finite-sample performance by our novel batch and online algorithms compared to existing approximate minimax algorithms.

2. Normalized Maximum Likelihood and Asymptotic Minimaxity

Consider a sequence $x^n = (x_1, \dots, x_n)$ and a parametric model

$$p(x^n|\theta) = \prod_{i=1}^n p(x_i|\theta),$$

where $\theta = (\theta_1, \dots, \theta_d)$ is a d -dimensional parameter. We focus on the case where each x_i is one of a finite alphabet of symbols and the maximum likelihood estimator

$$\hat{\theta}(x^n) = \operatorname{argmax}_{\theta} \ln p(x^n | \theta)$$

can be computed.

The optimal solution to the minimax problem,

$$\min_{\bar{p}} \max_{x^n} \ln \frac{p(x^n | \hat{\theta}(x^n))}{\bar{p}(x^n)}$$

is given by

$$p_{\text{NML}}^{(n)}(x^n) = \frac{p(x^n | \hat{\theta}(x^n))}{C_n},$$

where $C_n = \sum_{x^n} p(x^n | \hat{\theta}(x^n))$ and is called the normalized maximum likelihood (NML) distribution (Shtarkov, 1987). The minimax regret is given by $\ln C_n$ for all x^n . We mention that in addition to coding and prediction, the code length $-\ln p_{\text{NML}}^{(n)}(x^n)$ has been used as a model-selection criterion (Rissanen, 1996); see also Grünwald (2007); Silander et al. (2010) and references therein.

Since the normalizing constant C_n is computationally intractable in most models, we consider approximating the minimax optimal NML model by another model $g(x^n)$ and focus on *asymptotic* minimax optimality of g , which is defined by

$$\max_{x^n} \ln \frac{p(x^n | \hat{\theta}(x^n))}{g(x^n)} \leq \ln C_n + o(1), \quad (1)$$

where $o(1)$ is a term converging to zero as $n \rightarrow \infty$.

Under the following assumption, we can show (Thm. 1 below) that the model g must be dependent on the sample size n to achieve the asymptotic minimax optimality in a slightly stronger sense, as characterized in the theorem.

Assumption 1 *Suppose that for \tilde{n} satisfying $\tilde{n} \rightarrow \infty$ and $\frac{\tilde{n}}{n} \rightarrow 0$ as $n \rightarrow \infty$ (e.g. $\tilde{n} = \sqrt{n}$), there exist a sequence $x^{\tilde{n}}$ and a unique constant $M > 0$ such that*

$$\ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{\sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)} \rightarrow M \quad (n \rightarrow \infty), \quad (2)$$

where $\sum_{x_{\tilde{n}+1}^n} = \sum_{x_{\tilde{n}+1}} \cdots \sum_{x_n}$ denotes the marginalization over $x_{\tilde{n}+1}, \dots, x_n$.

Assumption 1 means that the NML model changes over the sample size n , the amount of which is characterized by M . The following theorem proves that under this assumption, the asymptotic minimaxity is never achieved simultaneously for the sample sizes \tilde{n} and n by an online strategy g that is independent of n .

Theorem 1 *If the model g is independent of the sample size n and satisfies $\sum_{x_{\tilde{n}+1}^n} g(x^n) = g(x^{\tilde{n}})$, then it never satisfies*

$$\ln C_n - \underline{M} + o(1) \leq \ln \frac{p(x^n|\hat{\theta}(x^n))}{g(x^n)} \leq \ln C_n + o(1), \quad (3)$$

for all x^n and any $\underline{M} < M$, where M is the constant appearing in Assumption 1 and $o(1)$ is a term converging to zero uniformly on x^n as $n \rightarrow \infty$.

The proof is given in Appendix A.

Note that the condition in Eq. (3) is stronger than the usual asymptotic minimax optimality in Eq. (1), where only the right inequality in Eq. (3) is required. Intuitively, our stronger notion of asymptotic minimaxity requires not only that for all sequences, the regret of the model g is asymptotically at most the minimax value, but also that for *no* sequence, the regret is asymptotically *less* than the minimax value by a margin characterized by \underline{M} . Note that non-asymptotically (without the $o(1)$ terms), the corresponding strong and weak minimax notions are equivalent since reducing the code length for one sequence (compared to the NML model), necessarily increases the code length for at least one other sequence.

When we take g as a Bayes mixture,

$$g(x^n) = \int p(x^n|\theta)q(\theta)d\theta,$$

$\sum_{x_{\tilde{n}+1}^n} g(x^n) = g(x^{\tilde{n}})$ holds if the prior distribution $q(\theta)$ does not depend on n . On the contrary, if $q(\theta)$ depends on n , it is possible that the Bayes mixture achieves Eq. (3) for all x^n . In fact, for the multinomial model (with m categories), the Dirichlet prior $\text{Dir}(\alpha_n, \dots, \alpha_n)$ with $\alpha_n = \frac{1}{2} - \frac{\ln 2}{2} \frac{1}{\ln n}$ provides an example of such a case as will be proven in Sect. 3.2. Section 3.1 demonstrates that the sequence of all 1s (or all 2s, 3s, etc.) gives $M = \frac{m-1}{2} \ln 2$ in the multinomial model.

3. Asymptotic Minimality in Multinomial Model

Hereafter, we focus on the multinomial model with $x \in \{1, 2, \dots, m\}$,

$$p(x|\theta) = \theta_x, \quad \sum_{j=1}^m \theta_j = 1.$$

Although a linear-time (in n) algorithm has been obtained for computing the NML distribution of this model (Kontkanen and Myllymäki, 2007), we examine asymptotic minimality of other strategies for this model whose theoretical properties have been studied in depth (Xie and Barron, 2000).

For the multinomial model, the Dirichlet distribution is a conjugate prior, taking the form

$$q(\theta) = \frac{\Gamma(m\alpha)}{\Gamma(\alpha)^m} \prod_{j=1}^m \theta_j^{\alpha-1},$$

where $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$ is the gamma function and $\alpha > 0$ is a hyperparameter. The Bayes mixture is obtained as follows,

$$\begin{aligned} p_{B,\alpha}(x^n) &= \int \prod_{i=1}^n p(x_i|\theta) q(\theta) d\theta \\ &= \frac{\Gamma(m\alpha) \prod_{j=1}^m \Gamma(n_j + \alpha)}{\Gamma(\alpha)^m \Gamma(n + m\alpha)}, \end{aligned} \quad (4)$$

where n_j is the number of j s in x^n . The minimax regret is asymptotically given by (Xie and Barron, 2000)

$$\ln C_n = \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \frac{\Gamma(1/2)^m}{\Gamma(m/2)} + o(1). \quad (5)$$

In the following two subsections, we evaluate the constant M of Assumption 1 and derive an asymptotically minimax optimal hyperparameter α . We use the following lemma in this section. The proof is in Appendix B.

Lemma 2 *Let*

$$f(x) = \ln \Gamma\left(x + \frac{1}{2}\right) - x \ln x + x - \frac{1}{2} \ln \pi.$$

Then for $x \geq 0$,

$$0 \leq f(x) < \frac{\ln 2}{2} \quad (6)$$

and $\lim_{x \rightarrow \infty} f(x) = \frac{\ln 2}{2}$.

3.1. Change of NML Model

Let l_j be the number of j s in $x^{\tilde{n}}$ ($0 \leq l_j \leq \tilde{n}$, $\sum_{j=1}^m l_j = \tilde{n}$). It follows that

$$\ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{\sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)} = \ln \frac{\prod_{j=1}^m \binom{l_j}{\tilde{n}}^{l_j}}{\sum_{n_j \geq l_j} \binom{n-\tilde{n}}{n_j-l_j} \prod_{j=1}^m \binom{n_j}{n}^{n_j}} + \ln \frac{C_n}{C_{\tilde{n}}}, \quad (7)$$

where $\binom{n-\tilde{n}}{n_j-l_j} \equiv \binom{n-\tilde{n}}{n_1-l_1, \dots, n_m-l_m}$ is the multinomial coefficient and $\sum_{n_j \geq l_j}$ denotes the summation over n_j s satisfying $n_1 + \dots + n_m = n$ and $n_j \geq l_j$ for $j = 1, 2, \dots, m$. The following lemma evaluates

$$C_{n|x^{\tilde{n}}} \equiv \sum_{n_j \geq l_j} \binom{n-\tilde{n}}{n_j-l_j} \prod_{j=1}^m \binom{n_j}{n}^{n_j}$$

in Eq. (7). The proof is in Appendix C.¹

1. For the Fisher information matrix $I(\theta)$ whose ij th element is given by $(I(\theta))_{ij} = -\sum_x p(x|\theta) \frac{\partial^2 \ln p(x|\theta)}{\partial \theta_i \partial \theta_j} = \delta_{i,j}/\theta_j$, the constant $\tilde{C}_{1/2}$ coincides with $\int \sqrt{|I(\theta)|} \prod_{j=1}^m \theta^{l_j} d\theta$. This proves that the asymptotic expression of the regret of the conditional NML (Grünwald, 2007, Eq. (11.47), p.323) is valid for the multinomial model with the full parameter set rather than the restricted parameter set discussed in Grünwald (2007).

Lemma 3 $C_{n|x^{\tilde{n}}}$ is asymptotically evaluated as

$$\ln C_{n|x^{\tilde{n}}} = \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \tilde{C}_{\frac{1}{2}} + o(1), \quad (8)$$

where \tilde{C}_{α} is defined for $\alpha > 0$ and $\{l_j\}_{j=1}^m$ as

$$\tilde{C}_{\alpha} = \frac{\prod_{j=1}^m \Gamma(l_j + \alpha)}{\Gamma(\tilde{n} + m\alpha)}. \quad (9)$$

Substituting Eq. (8) and Eq. (5) into Eq. (7), we have

$$\ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{p_{\text{NML}}^{(n)}(x^{\tilde{n}})} = -\frac{m-1}{2} \ln \frac{\tilde{n}}{2\pi} + \sum_{j=1}^m l_j \ln \frac{l_j}{\tilde{n}} - \ln \frac{\prod_{j=1}^m \Gamma(l_j + 1/2)}{\Gamma(\tilde{n} + m/2)} + o(1),$$

where $p_{\text{NML}}^{(n)}(x^{\tilde{n}}) = \sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)$. Applying Stirling's formula to $\ln \Gamma(\tilde{n} + m/2)$ expresses the right hand side as

$$\sum_{j=1}^m \left\{ l_j \ln l_j - \ln \Gamma \left(l_j + \frac{1}{2} \right) - l_j + \frac{1}{2} \ln 2\pi \right\} + o(1).$$

Taking $l_1 = \tilde{n}$, $l_j = 0$ for $j = 2, \dots, m$, from Lemma 2, we have $\ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{p_{\text{NML}}^{(n)}(x^{\tilde{n}})} = \frac{m-1}{2} \ln 2 + o(1)$, that is, Assumption 1 holds with $M = (m-1) \ln 2/2$.

3.2. Optimal Hyperparameter and its Asymptotic Minimacity

We examine the asymptotic minimacity of the Bayes mixture in Eq. (4). More specifically, we investigate the minimax optimal hyperparameter

$$\operatorname{argmin}_{\alpha} \max_{x^n} \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,\alpha}(x^n)} \quad (10)$$

and show that it is asymptotically approximated by

$$\alpha_n = \frac{1}{2} - \frac{\ln 2}{2} \frac{1}{\ln n}. \quad (11)$$

We assume that the maximum regret is attained by both x^n consisting of a single symbol repeated n times as well as x^n with a uniform number n/m of each symbol j .² Let the regrets of these two cases be equal,

$$\Gamma(\alpha)^{m-1} \Gamma(n + \alpha) = \Gamma(n/m + \alpha)^m m^n.$$

Taking logarithms, using Stirling's formula and ignoring diminishing terms, we have

$$\begin{aligned} & (m-1) \left(\alpha - \frac{1}{2} \right) \ln n - (m-1) \ln \Gamma(\alpha) \\ & - m \left(\alpha - \frac{1}{2} \right) \ln m + (m-1) \frac{\ln 2\pi}{2} = 0. \end{aligned} \quad (12)$$

2. This assumption is implied from the proof of Thm. 4 (see the note after the proof).

This implies that the optimal α is asymptotically given by

$$\alpha_n \simeq \frac{1}{2} - \frac{a}{\ln n}, \quad (13)$$

for some constant a . Substituting this back into Eq. (12) and solving it for a , we obtain Eq. (11).

We numerically calculated the optimal hyperparameter defined by Eq. (10) for the binomial model ($m = 2$). Figure 1 shows the optimal α obtained numerically and its asymptotic approximation in Eq. (11). We see that the optimal hyperparameter is well approximated by α_n in Eq. (11) for large n . Note here the slow convergence speed, $O(1/\ln n)$ to the asymptotic value, $1/2$.

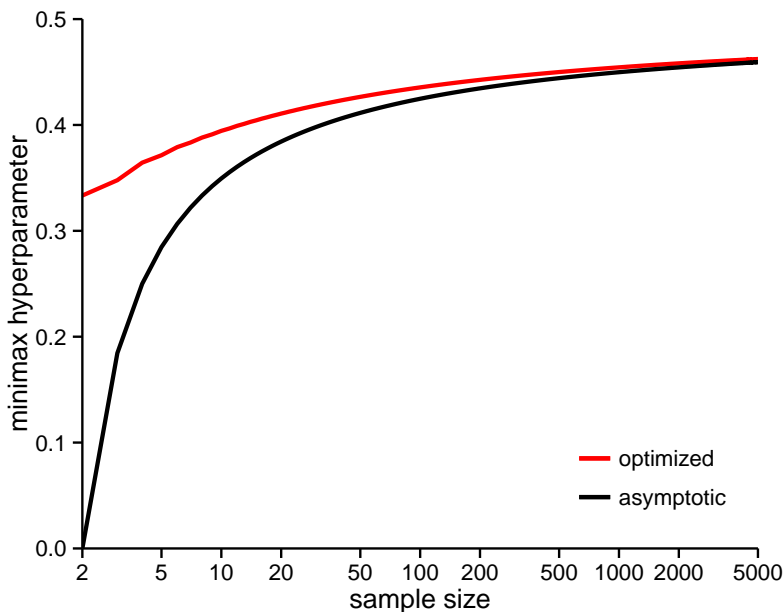


Figure 1: Minimax optimal hyperparameter α for sample size n

The next theorem shows the asymptotic minimaxity of α_n . It also shows that the lower bound in Eq. (3) is almost attainable for the multinomial model. We will examine the regret of α_n numerically in Sect. 4.1.

Theorem 4 *The Bayes mixture defined by the prior $\text{Dir}(\alpha_n, \dots, \alpha_n)$ is asymptotically minimax and satisfies*

$$\ln C_n - M + o(1) \leq \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B, \alpha_n}(x^n)} \leq \ln C_n + o(1), \quad (14)$$

for all x^n , where $M = (m - 1) \ln 2/2$.

The proof is given in Appendix D.

4. Numerical Results

In this section, we numerically calculate the maximum regrets of several methods in the binomial model ($m = 2$). The following two subsections respectively examine batch algorithms based on Bayes mixtures with prior distributions dependent on n and last-step minimax algorithms, which are online.

4.1. Optimal Conjugate Prior and Modified Jeffreys Prior

We calculated the maximum regrets of the Bayes mixtures in Eq. (4) with the hyperparameter optimized by the golden section search and with its asymptotic approximation in Eq. (11). We also investigated the maximum regrets of Xie and Barron’s modified Jeffreys prior which is proved to be asymptotically minimax (Xie and Barron, 2000). The modified Jeffreys prior is defined by

$$q_{\text{MJ}}^{(n)}(\theta) = \frac{\epsilon_n}{2} \left\{ \delta\left(\theta - \frac{1}{n}\right) + \delta\left(\theta - 1 + \frac{1}{n}\right) \right\} + (1 - \epsilon_n)b_{1/2}(\theta),$$

where δ is the Dirac’s delta function and $b_{1/2}(\theta)$ is the density function of the beta distribution with hyperparameters $1/2$, $\text{Beta}(1/2, 1/2)$, which is the Jeffreys prior for the Bernoulli model. We set $\epsilon_n = n^{-1/8}$ as proposed in Xie and Barron (2000) and also optimized ϵ_n by the golden section search so that the maximum regret

$$\max_{x^n} \ln \frac{p(x^n | \hat{\theta}(x^n))}{\int p(x^n | \theta) q_{\text{MJ}}^{(n)}(\theta) d\theta}$$

is minimized.

Figure 2(a) shows the maximum regrets of these Bayes mixtures: asymptotic and optimized Beta refer to mixtures with Beta priors (Sect. 3.2), and modified Jeffreys methods refer to mixtures with a modified Jeffreys prior as discussed above. Also included for comparison is the maximum regret of the Jeffreys mixture (Krichevsky and Trofimov, 1981), which is known not to be asymptotically minimax. To better show the differences, the regret of the NML model, $\ln C_n$, is subtracted from the maximum regret of each model.

We see that the maximum regrets of these models, except the one based on Jeffreys prior, decrease toward zero as n grows as implied by their asymptotic minimaxity. The modified Jeffreys prior with the optimized weight performs best of these strategies for this range of the sample size while that with the unoptimized weight performs much worse. Note here that we have the explicit form of the asymptotically minimax hyperparameter in Eq. (11) whereas the optimal weight for the modified Jeffreys prior is not known analytically. Note also that unlike the NML, Bayes mixtures can be computed in a sequential manner with respect to (x_1, \dots, x_n) even if the prior depends on n . The time complexity for online prediction will be discussed in Sect. 4.3.

The differences in the maximum regrets under the binomial model in Fig. 2 are small (less than 1 nat). However, they may be important even from a practical point of view. For instance, it has been empirically observed that the slightest differences in the Dirichlet hyperparameter can be significant in Bayesian network structure learning (Silander et al., 2007). Furthermore, the differences are likely to be greater under multinomial ($m > 2$) and other kinds of models.

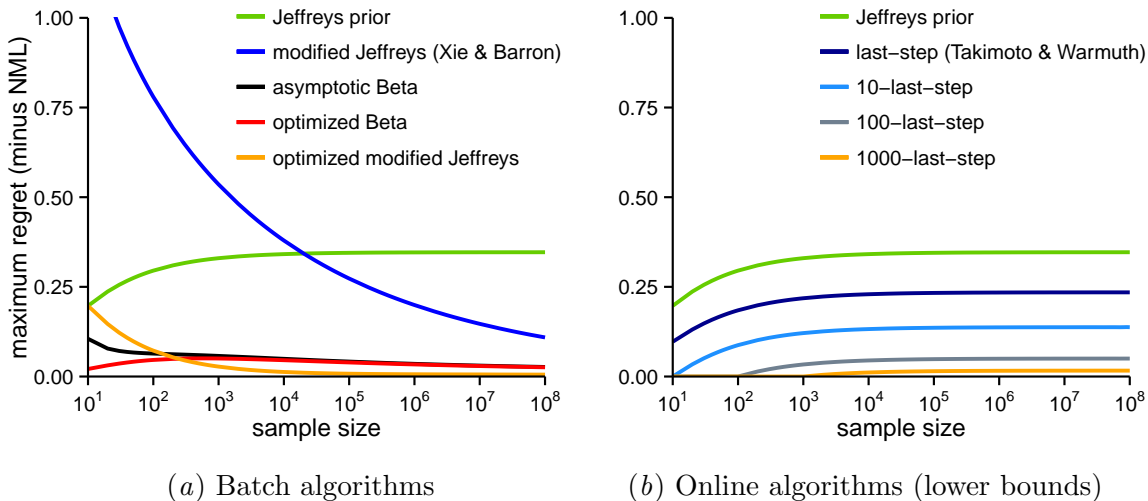


Figure 2: Maximum regret for sample size n . The regret of the NML model, $\ln C_n$, is subtracted from the maximum regret of each strategy. The first two algorithms in each panel are from earlier work, while the remaining ones are novel.

In Fig. 3, we show the *minimum* regrets of the Bayes mixtures with the optimal hyperparameter and that with its asymptotic approximation α_n to confirm the lower bound in Eq. (14). As we proved in Thm. 4, the regret (minus $\ln C_n$) is greater than the lower bound, $-M = -\ln 2/2$. The minimum regrets of the modified Jeffreys mixtures (not shown) were much smaller than the lower bound, $-\ln 2/2$. This implies that the modified Jeffreys mixture provides an example of an asymptotically minimax strategy but not in the sense of Eq. (3).

4.2. Last-Step Minimax Algorithms

The last-step minimax algorithm is an online prediction algorithm that is equivalent to the so called sequential normalized maximum likelihood method in the case of the multinomial model (Rissanen and Roos, 2007; Takimoto and Warmuth, 2000). A straightforward generalization, which we call the k -last-step minimax algorithm, normalizes $p(x^t|\hat{\theta}(x^t))$ over the last $k \geq 1$ steps to calculate the conditional distribution of $x_{t-k+1}^t = \{x_{t-k+1}, \dots, x_t\}$,

$$p_{\text{kLS}}(x_{t-k+1}^t|x^{t-k}) = \frac{p(x^t|\hat{\theta}(x^t))}{L_{t,k}},$$

where $L_{t,k} = \sum_{x_{t-k+1}^t} p(x^t|\hat{\theta}(x^t))$. Although this generalization was mentioned in Takimoto and Warmuth (2000), it was left as an open problem to examine how k affects the regret of the algorithm.

It is suggested from Thm. 1 that k -last-step minimax algorithm with k independent of n is not asymptotically minimax (although the theorem does not rigorously exclude

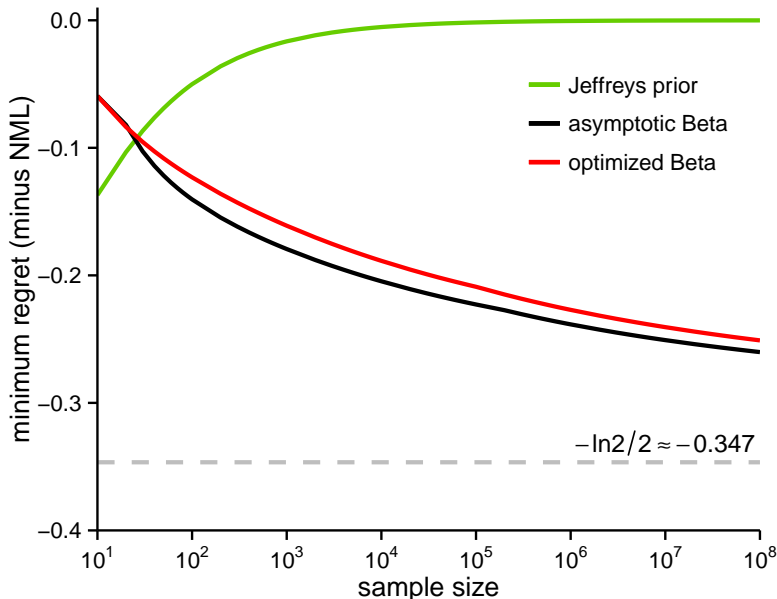


Figure 3: Minimum regret (minus $\ln C_n$) for sample size n .

that possibility because of the left inequality in Eq. (3)). We numerically calculated the regret of the k -last-step minimax algorithm with $k = 1, 10, 100$ and 1000 for the sequence $x^n = 1010101010 \dots$ since it is infeasible to evaluate the maximum regret for large n . The regret for this particular sequence provides a lower bound for the maximum regret. Figure 2(b) shows the regret as a function of n together with the maximum regret of the Jeffreys mixture. The theoretical asymptotic regret for the Jeffreys mixture is $\frac{\ln 2}{2} \approx 0.34$ (Krichevsky and Trofimov, 1981), and the asymptotic bound for the 1-last-step minimax algorithm is slightly better, $\frac{1}{2} (1 - \ln \frac{\pi}{2}) \approx 0.27$ (Takimoto and Warmuth, 2000). We can see that although the regret decreases as k grows, it still increases as n grows and does not converge to that of the NML (zero in the figure).

4.3. Computational Complexity

Although the NML distribution of the multinomial model is computed in linear time in n (Kontkanen and Myllymäki, 2007), the algorithm only provides the total code-length (or probability) of any complete sequence x^n . For prediction purposes in online learning scenarios, NML requires to compute the predictive probabilities $p_{\text{NML}}^{(n)}(x_t|x^{t-1})$ by summing over all continuations of x^t . Computing all the predictive probabilities up to n takes the time complexity of $O(m^n)$. For all the other algorithms except NML, the complexity is $O(n)$. More specifically, for Bayes mixtures, the complexity is $O(mn)$ and for k -laststep minimax algorithms, the complexity is $O(m^k n)$.

5. Discussion & Conclusion

In this paper, we proved that the knowledge of the sample size n is required for a strategy to be asymptotically minimax in the sense of Eq. (3). Bartlett et al. (2013) proved, as a corollary to their main result, that NML is sample-size dependent in the general exponential family. We have not observed any asymptotically minimax strategy independent of n . This suggests that the lower bound in Eq. (3) may be removed from the condition and the asymptotic minimaxity in the usual sense may be characterized by the dependency on n ; in other words, no online strategy can be asymptotically minimax.

For the multinomial model, Thm. 4 shows that a simple dependency on n is sufficient to provide an accurate approximation. In practice, our numerical experiments suggest the superiority of a number of novel algorithms, whose performance is very near that of the NML, both of the batch (Fig. 2(a)) as well as online (Fig. 2(b)) type.

The Dirichlet prior $\text{Dir}(\alpha, \dots, \alpha)$ has yielded related estimators for the multinomial model. The uniform prior, $\alpha = 1$, yields the Laplace estimator and the Jeffreys prior, $\alpha = 0.5$, yields the Krichevsky-Trofimov estimator (Krichevsky and Trofimov, 1981). Krichevsky (1998) showed that $\alpha = 0.50922\dots$ is optimal when the goal is to minimize the worst-case expected redundancy in predicting the $(n + 1)$ st symbol after a sequence of n symbols. Komaki (2012) studied a similar one-step ahead prediction where the boundary of the parameter space is treated in a certain way and obtained another prior, $\alpha = 1 + \sqrt{6}$. Tjalkens et al. (1993) and Hutter (2013) propose estimators whose regrets have the leading term $\frac{\tilde{m}-1}{2} \ln n$, where \tilde{m} is the number of different symbols that appear in x^n . These estimators are obtained by the prior depending on \tilde{m} and n and are designed for the case of a large alphabet size where $n \gg m$ cannot be expected. For a fixed n , however, the NML minimizes the worst-case regret over all possible sequences. The minimax optimal estimator among the Dirichlet-multinomial model is approximated by $\alpha_n = \frac{1}{2} - \frac{\ln 2}{2 \ln n}$, which is asymptotically optimal up to $o(1)$ terms as was proved in Sect. 3.2.

Future directions include verifying our conjecture about non-achievability of minimax regret by online strategies, developing approximation schemes of the NML distribution for more complex models, and their applications in prediction, data compression, and model selection.

Acknowledgments

This work was supported in part by JSPS KAKENHI Grant Numbers 23700175, 25280083, 25120014, by the Academy of Finland (project Prime and the COIN Center of Excellence) and by the Finnish Funding Agency for Technology and Innovation (project D2I). Part of this work was carried out while the first author was visiting HIIT.

References

- K. S. Azoury and M. K. Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. *Machine Learning*, 43(3):211–246, 2001.
- P. Bartlett, P. Grünwald, P. Harremoës, F. Hedayati, and W. Kotłowski. Horizon-independent optimal prediction with log-loss in exponential families. In *JMLR: Workshop*

- and *Conference Proceedings: 26th Annual Conference on Learning Theory*, volume 30, pages 639–661, 2013.
- N. Cesa-Bianchi and G. Lugosi. Worst-case bounds for the logarithmic loss of predictors. *Machine Learning*, 43(3):247–264, 2001.
- Y. Freund. Predicting a binary sequence almost as well as the optimal biased coin. In *Proc. of Computational Learning Theory (COLT' 96)*, pages 89–98, 1996.
- P. D. Grünwald. *The Minimum Description Length Principle*. The MIT Press, 2007.
- M. Hutter. Sparse adaptive Dirichlet-multinomial-like processes. In *JMLR: Workshop and Conference Proceedings: 26th Annual Conference on Learning Theory*, volume 30, pages 432–459, 2013.
- F. Komaki. Asymptotically minimax Bayesian predictive densities for multinomial models. *Electronic Journal of Statistics*, 6:934–957, 2012.
- P. Kontkanen and P. Myllymäki. A linear-time algorithm for computing the multinomial stochastic complexity. *Information Processing Letters*, 103(6):227–233, 2007.
- W. Kotłowski and P. D. Grünwald. Maximum likelihood vs. sequential normalized maximum likelihood in on-line density estimation. In *JMLR: Workshop and Conference Proceedings: 24th Annual Conference on Learning Theory*, volume 19, pages 457–476, 2011.
- R. E. Krichevsky. Laplace’s law of succession and universal encoding. *IEEE Trans. Information Theory*, 44(1):296–303, 1998.
- R. E. Krichevsky and V. K. Trofimov. The performance of universal encoding. *IEEE Trans. Information Theory*, IT-27(2):199–207, 1981.
- M. Merkle. Conditions for convexity of a derivative and some applications to the Gamma function. *Aequationes Mathematicae*, 55:273–280, 1998.
- J. Rissanen. Fisher information and stochastic complexity. *IEEE Trans. Information Theory*, IT-42(1):40–47, 1996.
- J. Rissanen and T. Roos. Conditional NML universal models. In *Proc. 2007 Information Theory and Applications Workshop (ITA-2007)*, pages 337–341. IEEE Press, 2007.
- Y. M. Shtarkov. Universal sequential coding of single messages. *Problems of Information Transmission*, 23(3):175–186, 1987.
- T. Silander, T. Roos, and P. Myllymäki. Learning locally minimax optimal Bayesian networks. *International Journal of Approximate Reasoning*, 51(5):544–557, 2010.
- Tomi Silander, Petri Kontkanen, and Petri Myllymäki. On sensitivity of the MAP Bayesian network structure to the equivalent sample size parameter. In *UAI*, pages 360–367, 2007.
- J. Takeuchi and A. R. Barron. Asymptotically minimax regret for exponential families. In *Proc. of the 20th Symposium on Information Theory and its Applications (SITA '97)*, pages 665–668, 1997.

- E. Takimoto and M. K. Warmuth. The last-step minimax algorithm. In *Algorithmic Learning Theory, Lecture Notes in Computer Science*, volume 1968, pages 279–290, 2000.
- T. J. Tjalkens, Y. M. Shtarkov, and F. M. J. Willems. Sequential weighting algorithms for multi-alphabet sources. In *Proc. 6th Joint Swedish-Russian Intl. Workshop on Information Theory*, pages 22–27, 1993.
- Q. Xie and A. R. Barron. Asymptotic minimax regret for data compression, gambling, and prediction. *IEEE Trans. Information Theory*, 46(2):431–445, 2000.

Appendix A. Proof of Theorem 1

Proof Under Assumption 1, we suppose Eq. (3) holds for all sufficiently large n and derive contradiction. The inequalities in Eq. (3) are equivalent to

$$-\underline{M} + o(1) \leq \ln \frac{p_{\text{NML}}^{(n)}(x^n)}{g(x^n)} \leq o(1).$$

This implies that

$$\begin{aligned} g(x^{\tilde{n}}) &= \sum_{x_{\tilde{n}+1}^n} g(x^n) = \sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n) e^{-\ln \frac{p_{\text{NML}}^{(n)}(x^n)}{g(x^n)}} \\ &\leq e^{\underline{M}+o(1)} \sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n) \end{aligned} \quad (15)$$

for all $x^{\tilde{n}}$. Then we have

$$\begin{aligned} \max_{x^{\tilde{n}}} \ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{g(x^{\tilde{n}})} &= \max_{x^{\tilde{n}}} \left\{ \ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{\sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)} + \ln \frac{\sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)}{g(x^{\tilde{n}})} \right\} \\ &\geq \max_{x^{\tilde{n}}} \left\{ \ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{\sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)} \right\} - \underline{M} + o(1) \\ &\geq \epsilon + o(1), \end{aligned}$$

where $\epsilon = M - \underline{M} > 0$. The first inequality follows from Eq. (15) and the second inequality follows from Assumption 1, which implies $\max_{x^{\tilde{n}}} \ln \frac{p_{\text{NML}}^{(\tilde{n})}(x^{\tilde{n}})}{\sum_{x_{\tilde{n}+1}^n} p_{\text{NML}}^{(n)}(x^n)} \geq M + o(1)$. The above inequality contradicts the asymptotic minimax optimality in Eq. (3) with n replaced by \tilde{n} . ■

Appendix B. Proof of Lemma 2

Proof The function f is non-decreasing since $f'(x) = \psi(x + 1/2) - \ln x \geq 0$ where $\psi(x) = (\ln \Gamma(x))'$ is the psi function (Merkle, 1998). $\lim_{x \rightarrow \infty} f(x) = \frac{\ln 2}{2}$ is derived from Stirling's

formula,

$$\ln \Gamma(x) = \left(x - \frac{1}{2}\right) \ln x - x + \frac{1}{2} \ln(2\pi) + O\left(\frac{1}{x}\right).$$

It immediately follows from $f(0) = 0$ and this limit that $0 \leq f(x) < \frac{\ln 2}{2}$ for $x \geq 0$. \blacksquare

Appendix C. Proof of Lemma 3

Proof In order to prove Lemma 3, we modify and extend Xie and Barron's proof in Xie and Barron (2000) for the asymptotic evaluation of $\ln C_n = \ln \sum_{x^n} p(x^n | \hat{\theta}(x^n))$ given by Eq. (5) to that of $\ln C_{n|x^{\tilde{n}}} = \ln \sum_{x_{\tilde{n}+1}^n} p(x^n | \hat{\theta}(x^n))$, which is conditioned on the first \tilde{n} samples, $x^{\tilde{n}}$. More specifically, we will prove the following inequalities. Here, $p_{B,w}$ denotes the Bayes mixture defined by the prior $w(\theta)$, $p_{B,1/2}$ and p_{B,α_n} are those with the Dirichlet priors, $\text{Dir}(1/2, \dots, 1/2)$ (Jeffreys mixture) and $\text{Dir}(\alpha_n, \dots, \alpha_n)$ where $\alpha_n = \frac{1}{2} - \frac{\ln 2}{2} \frac{1}{\ln n}$ respectively.

$$\begin{aligned} \frac{m-1}{2} \ln \frac{n}{2\pi} + \tilde{C}_{\frac{1}{2}} + o(1) &\leq \sum_{x_{\tilde{n}+1}^n} p_{B,1/2}(x_{\tilde{n}+1}^n | x^{\tilde{n}}) \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,1/2}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} & (16) \\ &\leq \max_w \sum_{x_{\tilde{n}+1}^n} p_{B,w}(x_{\tilde{n}+1}^n | x^{\tilde{n}}) \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,w}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} \\ &= \max_w \min_{\bar{p}} \sum_{x_{\tilde{n}+1}^n} p_{B,w}(x_{\tilde{n}+1}^n | x^{\tilde{n}}) \ln \frac{p(x^n | \hat{\theta}(x^n))}{\bar{p}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} \\ &\leq \min_{\bar{p}} \max_{x_{\tilde{n}+1}^n} \ln \frac{p(x^n | \hat{\theta}(x^n))}{\bar{p}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} \\ &= \ln \sum_{x_{\tilde{n}+1}^n} p(x^n | \hat{\theta}(x^n)) = \ln C_{n|x^{\tilde{n}}} \\ &\leq \max_{x_{\tilde{n}+1}^n} \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,\alpha_n}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} \\ &\leq \frac{m-1}{2} \ln \frac{n}{2\pi} + \tilde{C}_{\frac{1}{2}} + o(1), & (17) \end{aligned}$$

where the first equality follows from Gibbs' inequality, and the second equality as well as the second to last inequality follow from the minimax optimality of NML (Shtarkov, 1987). Let us move on to the proof of inequalities (16) and (17). The rest of the inequalities follow from the definitions and from the fact that maximin is no greater than minimax. To derive both inequalities, we evaluate $\ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,\alpha}(x_{\tilde{n}+1}^n | x^{\tilde{n}})}$ for the Bayes mixture with the prior

$\text{Dir}(\alpha, \dots, \alpha)$ asymptotically. It follows that

$$\begin{aligned}
 \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,\alpha}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} &= \ln \frac{\prod_{j=1}^m \left(\frac{n_j}{n}\right)^{n_j}}{\frac{\Gamma(\tilde{n}+m\alpha)}{\Gamma(n+m\alpha)} \prod_{j=1}^m \frac{\Gamma(n_j+\alpha)}{\Gamma(l_j+\alpha)}} \\
 &= \sum_{j=1}^m n_j \ln n_j - n \ln n - \sum_{j=1}^m \ln \Gamma(n_j + \alpha) + \ln \Gamma(n + m\alpha) + \ln \tilde{C}_\alpha \\
 &= \sum_{j=1}^m \left\{ n_j \ln n_j - n_j - \ln \Gamma(n_j + \alpha) + \frac{1}{2} \ln(2\pi) \right\} \\
 &\quad + \left(m\alpha - \frac{1}{2}\right) \ln n - (m-1) \frac{1}{2} \ln(2\pi) + \ln \tilde{C}_\alpha + o(1), \tag{18}
 \end{aligned}$$

where \tilde{C}_α is defined in Eq. (9) and we applied Stirling's formula to $\ln \Gamma(n + m\alpha)$.

Substituting $\alpha = 1/2$ into Eq. (18), we have

$$\ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,1/2}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} = \sum_{j=1}^m \left(c_{n_j} + \frac{\ln 2}{2} \right) + \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \tilde{C}_{1/2} + o(1),$$

where

$$c_k = k \ln k - k - \ln \Gamma(k + 1/2) + \frac{1}{2} \ln \pi, \tag{19}$$

for $k \geq 0$. Since from Lemma 2, $-\frac{\ln 2}{2} < c_k$,

$$\ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,1/2}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} > \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \tilde{C}_{1/2} + o(1),$$

holds for all x^n , which proves the inequality (16).

Substituting $\alpha = \alpha_n = \frac{1}{2} - \frac{\ln 2}{2} \frac{1}{\ln n}$ into Eq. (18), we have

$$\begin{aligned}
 \ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,\alpha_n}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} &= \sum_{j=1}^m \left\{ n_j \ln n_j - n_j - \ln \Gamma(n_j + \alpha_n) + \frac{1}{2} \ln \pi \right\} \\
 &\quad + \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \tilde{C}_{1/2} + o(1).
 \end{aligned}$$

Assuming that the first l n_j s ($j = 1, \dots, l$) are finite and the rest are large (tend to infinity as $n \rightarrow \infty$) and applying Stirling's formula to $\ln \Gamma(n_j + \alpha_n)$ ($j = l+1, \dots, m$), we have

$$\ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B,\alpha_n}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} = \sum_{j=1}^l c_{n_j} + \sum_{j=l+1}^m d_{n_j} + \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \tilde{C}_{1/2} + o(1), \tag{20}$$

where c_k is defined in Eq. (19) and

$$d_k = \frac{\ln 2}{2} \left(\frac{\ln k}{\ln n} - 1 \right)$$

for $1 < k \leq n$. Since $c_k \leq 0$ follows from Lemma 2 and $d_k \leq 0$, we obtain

$$\ln \frac{p(x^n | \hat{\theta}(x^n))}{p_{B, \alpha_n}(x_{\tilde{n}+1}^n | x^{\tilde{n}})} \leq \frac{m-1}{2} \ln \frac{n}{2\pi} + \ln \tilde{C}_{1/2} + o(1), \quad (21)$$

for all x^n , which proves the inequality (17). \blacksquare

Appendix D. Proof of Theorem 4

Proof The proof of Lemma 3 itself applies to the case where $\tilde{n} = 0$ and $l_j = 0$ for $j = 1, \dots, m$ as well. Since, in this case, $\tilde{C}_{1/2} = \ln \frac{\Gamma(1/2)^m}{\Gamma(m/2)}$, Eq. (21) in the proof gives the right inequality in Eq. (14).

Furthermore, in Eq. (20), we have

$$\sum_{j=1}^l c_{n_j} + \sum_{j=l+1}^m d_{n_j} > -(m-1) \frac{\ln 2}{2} + o(1). \quad (22)$$

This is because, from Lemma 2 and definition, $c_{n_j}, d_{n_j} > -\frac{\ln 2}{2}$ and for at least one of j , n_j is in the order of n since $\sum_{j=1}^n n_j = n$, which means that $d_{n_j} = o(1)$ for some j . Substituting Eq. (22) into Eq. (20), we obtain the left inequality in Eq. (14) with $M = (m-1) \ln 2/2$. \blacksquare

In Eq. (20), $c_{n_j} = 0$ holds only for $n_j = 0$ and $d_{n_j} = o(1)$ holds only when n_j is of the order of n , say, n/m . This means that the maximum regret is obtained at the boundary $n_j = 0$ or around $n_j \simeq n/m$.