

# A Lower Bound for Regret in Logistic Regression

Gil I. Shamir  
Google Inc.  
Pittsburgh, PA, USA  
Email: gshamir@ieee.org

Wojciech Szpankowski  
Purdue University  
West Lafayette, IN, USA  
Email: szpan@purdue.edu

**Abstract**—We study logistic regression with binary features in which the number (or degree) of occurring features determines the label probability. This model fits one of social networks, where the number of friends determines the likelihood of outcomes instead of the identity of the friends, or more generally, a graph model, where the degree of a node can determine its behavior. It includes the case in which weights can be viewed as i.i.d. (e.g., in Bayesian modeling). For such a model, we introduce the *maximal minimax regret* that we analyze using a unique combination of analytic combinatorics and information theory. More importantly, the resulting regret is a general lower bound for the pointwise regret of a general logistic regression over all algorithms (learning distributions). We show that the introduced worst case (maximum over feature sequences) maximal minimax regret grows asymptotically as  $(d/2)\log(T/d) + (d/2)\log(\pi/2) + O(d/\sqrt{T})$  for dimensionality  $d = o(\sqrt{T})$ , which is a lower bound for a regret of a general logistic regression. We extend our results to loss functions other than logistic loss and non-binary labels.

## I. INTRODUCTION

Logistic regression has recently received much attention in machine learning (see [22]) due to its importance in many practical systems and applications from category classification to risk assessment. It consists of a set of features, whose parameters represent their effect on some outcome. In a supervised online setup, a model is trained to learn these parameters from examples whose outcomes are already labeled. The training algorithm consumes data in rounds, where at each round  $t \in \{1, 2, \dots, T\}$ , it is allowed to predict the label based only on the labels it observed in the past  $t - 1$  rounds. The prediction algorithm incurs for each round some *loss* and updates its belief of the model parameters. The (pointwise) *regret* of an online algorithm is defined as the (excess) loss it incurs over some value of a constant *comparator* (weight vector) that is used for prediction of the complete sequence.

Here we introduce and analyze the *maximal minimax regret* that for a given feature sequence maximizes the pointwise regret over label sequences and minimizes over learned distributions. Such a minimax regret was analyzed in information theory in the context of universal compression as discussed in [26], [5] while (pointwise) regret for logistic regression has been studied in [9], [12], [16], [18], [15], [24].

In this paper, instead of learning an individual parameter per feature, we study the problem where the outcome or *label* observed for an example is a function of the number (degree) of nonzero binary features present in the example (see assumption (16) in Section III). This setup can characterize social networks, in which the likelihood of some outcome depends on the

number of friends a person has instead of on who those friends are. More generally, this setup can be representative of graphs, where an outcome in some node depends on the degree of the node, and not on which nodes are its neighbors (e.g., see graph structural compression in [2] and [17]). Given that the class of parameters for which we precisely compute the regret for this problem is a subset of the class of parameters in the standard logistic regression setup, this regret is a “universal” lower bound for a general pointwise regret.

We first focus on *binary* labels and consider the case of binary features  $\mathbf{x}_t$  (i.e., the examples or feature values vector at any round  $t$  is a sparse vector in  $\{0, 1\}^d$ , where  $d$  is the dimension – the number of features). For such a set up, we introduce the *maximal minimax regret* that we analyze using a unique combination of analytic combinatorics and information theory. In Theorem 1 we rigorously show that the introduced minimax regret of the degree problem serves as a general lower bound for a general pointwise regret over any algorithm/ learning distribution. With this in mind, we turn our attention to precise analysis of the introduced minimax regret. In Theorem 2 we show that the *average* maximal minimax regret grows asymptotically like  $\frac{d}{2}\log T + \frac{1}{2}\log(\alpha_1 \cdots \alpha_d) + \frac{d}{2}\log(\pi/2) + O(d/\sqrt{T})$  where  $T$  denotes the number of rounds, and  $\alpha_i$  is the fraction of feature vector  $\mathbf{x}_t$  with exactly  $j$  active features for  $t = 1, \dots, T$ . We also show that the worst case (maximum over all feature sequences) maximal minimax behaves asymptotically like  $\frac{d}{2}\log \frac{T}{d} + \frac{d}{2}\log(\pi/2) + O(d/\sqrt{T})$  for large  $T$ . In Theorem 4, we extend the result to non-binary labels of size  $m$  such that  $m^{3/2}d = o(\sqrt{T})$ . Furthermore, in Corollary 1 we summarize our findings regarding the lower bound on the standard regret. In Theorem 3, we briefly discuss Bayesian setting and present results for the pointwise and average regrets.

Studying this problem, we adapt techniques from the universal compression literature (see [26], [29], [30], [31]) and analytic combinatorics (see [7], [28]) that apply complex asymptotics to solve discrete problems. We first review various notions of regret and redundancy from information theory that we adopt to the performance evaluation of logistic regression. Here, we assume that  $d = 1$  and alphabet is of size  $m$ . The *pointwise redundancy*  $R_T(P; y^T)$  and the *average redundancy*  $\bar{R}_T(P)$  for a given source  $P$  and source (label) sequence  $y^T = (y_1, \dots, y_T)$  of length  $T$  over alphabet  $\mathcal{A}$  of size  $m$  are

defined as

$$\begin{aligned} R_T(P; y^T) &= L(y^T) + \log P(y^T), \\ \bar{R}_T(P) &= \mathbf{E}[L(Y^T)] - H_T(P), \end{aligned}$$

where  $H_T(P)$  is the entropy for a block of length  $T$ ,  $\mathbf{E}$  denotes the expectation, and  $L(y^T)$  is the code length of some code  $L(\cdot)$  (both redundancy quantities are implicitly functions of the code). In online learning – and indeed in information theory – one ignores the integer nature of the length and replace it by  $L(y^T) = -\log Q(y^T)$  for some unknown distribution  $Q$  that approximates  $P$ . The definitions above imply a probabilistic setting, in which there is some source that generated the data. A non-probabilistic setting considers *individual sequences* (see, e.g., [26]), where we define the *maximal* or *worst case* redundancy as

$$R_T^*(Q, P) = \max_{y^T} [-\log Q(y^T) + \log P(y^T)].$$

In practice, one can only hope to have some knowledge about a family of sources  $\mathcal{S}$  that generates real data. Following Davisson [4], we define the average minimax redundancy  $\bar{R}_T(\mathcal{S})$  and the worst case (maximal) minimax redundancy  $R_T^*(\mathcal{S})$  for family  $\mathcal{S}$  as follows

$$\begin{aligned} \bar{R}_T(\mathcal{S}) &= \min_Q \sup_{P \in \mathcal{S}} \sum_{y^T} P(y^T) \log [P(y^T)/Q(y^T)], \\ R_T^*(\mathcal{S}) &= \min_Q \sup_{P \in \mathcal{S}} \max_{y^T} \{ \log [P(y^T)/Q(y^T)] \}. \end{aligned}$$

In words, we search for the best code or distribution  $Q$  for the worst source  $P$  on average and for the worst label sequence  $y^T$  for individual sequences.

There are other measures of optimality for coding, gambling, learning, and prediction that are used in universal modeling, learning, and coding. We refer here to minimax *regrets* defined as follows (cf. [5], [30], [31]):

$$\begin{aligned} \bar{r}_T(\mathcal{S}) &= \min_Q \sup_{P \in \mathcal{S}} \mathbf{E}_P [-\log Q(y^T) + \log \sup_{P \in \mathcal{S}} P(y^T)], \\ r_T^*(\mathcal{S}) &= \min_Q \max_{y^T} [-\log Q(y^T) + \log \sup_{P \in \mathcal{S}} P(y^T)], \end{aligned}$$

and to the maxmin regret

$$r_T(\mathcal{S}) = \sup_{P \in \mathcal{S}} \min_Q \mathbf{E} [-\log Q(y^T) + \log \sup_{P \in \mathcal{S}} P(y^T)].$$

We call  $\bar{r}_T(\mathcal{S})$  the *average* minimax regret,  $r_T^*(\mathcal{S})$  the worst case (maximal) minimax regret and  $r_T(\mathcal{S})$  the maxmin regret. Clearly,  $\bar{R}_T(\mathcal{S}) \leq \bar{r}_T(\mathcal{S})$ , and,  $r_T^*(\mathcal{S}) = R_T^*(\mathcal{S})$ .

In [5] it is also shown that if the maximum likelihood distribution belongs to the convex hull of  $\mathcal{S}$ , then  $R_T^*(\mathcal{S}) - \bar{R}_T(\mathcal{S}) = O(c_T(\mathcal{S}))$  where

$$c_T(\mathcal{S}) = \sum_{y^T} P(y^T) \log \frac{\sup_{P \in \mathcal{S}} P(y^T)}{P(y^T)}.$$

Furthermore, it is known [5], [19], [21], [23], [27], [30], [31] that for a large class of sources (up to Markovian but not for non-Markovian [8], [5]) the redundancy grows as  $\frac{m-1}{2} \log T$

when the alphabet size  $m$  is fixed and  $\frac{m-1}{2} \log(T/m)$  for  $m = o(T)$  (see [19], [23], [29]). In fact in [29] full asymptotic expansions were derived for all ranges of  $m$ .

Finally, we review the ML literature with respect to pointwise regret. To the best of our knowledge, in the ML literature, [16] was first to demonstrate results that suggest that pointwise regret for logistic regression grows like  $O(d \log T/d)$  for fixed dimension  $d$  and  $m = 2$ . This was further generalized in [9] to all  $m$ . The authors of [16] used *Bayesian model averaging*. The  $O(\log T)$  pointwise and individual sequence regret can be achieved for the single dimensional problem with gradient methods based approaches, as was demonstrated in [18]. The authors of [18] then posed the problem of what happens for larger dimensions. Subsequently, [9] demonstrated how to achieve regret bounds of  $O(d \log(T/d))$  with Bayesian model averaging. These results were strengthened in [24], which also provided matching lower bounds (see also [13]). We should point out that our results on minimax regret, which we show below, are not restricted to Bayesian modeling.

## II. PROBLEM FORMULATION AND NOTATION

We denote by  $\mathbf{x}_t = (x_{1,t}, \dots, x_{d,t})$  a  $d$ -dimensional *binary* feature vector. Throughout we assume that a feature is binary, that is,  $x_{i,t} \in \{0, 1\}$  for  $t = 1, \dots, T$ . Notice that  $\mathbf{x}^T$  is a  $T \times d$  matrix with  $\mathbf{x}_t = (x_{1,t}, \dots, x_{d,t})$  as a row. The label binary vector is denoted as  $y^T = (y_1, \dots, y_T)$  with  $y_t \in \{-1, 1\}$ . Finally,  $\mathbf{w}_t = (w_{1,t}, \dots, w_{d,t})$  is a  $d$ -dimensional vector of feature weights. In this paper, we do not address the method used to learn the weights (e.g., gradient method or Bayesian mixing). Instead, we assume that the weights are *exchangeable* (see assumption (D) in the next section) leading to our model in which the number of features determines the label probability. Notice that a practical learning algorithm will be agnostic to the vector  $\mathbf{w}$ , and instead will learn “features”, which are a function of the degree of  $\mathbf{x}_t$ , i.e., the number of nonzero components in  $\mathbf{x}_t$ .

The *logistic loss* of an algorithm that *plays*  $\mathbf{w}_t$  at round  $t$  is

$$L(y^T | \mathbf{x}^T, \mathbf{w}^T) := \sum_{t=1}^T \log [1 + \exp(-y_t \langle \mathbf{x}_t, \mathbf{w}_t \rangle)] \quad (1)$$

where  $\langle \mathbf{x}_t, \mathbf{w}_t \rangle = \sum_{i=1}^d x_{i,t} w_{i,t}$ . It is convenient to write  $\ell(y_t | \mathbf{x}_t, \mathbf{w}_t) := \log [1 + \exp(-y_t \langle \mathbf{x}_t, \mathbf{w}_t \rangle)]$ . Both  $\ell(y_t | \mathbf{x}_t, \mathbf{w}_t)$  and  $L(y^T | \mathbf{x}^T, \mathbf{w}^T)$  depend on  $\mathbf{x}_t$  and  $\mathbf{w}_t$  only through the product  $\langle \mathbf{x}_t, \mathbf{w}_t \rangle$ .

The probability of a label is given by

$$P(y_t | \mathbf{x}_t, \mathbf{w}_t) = \frac{1}{1 + \exp(-y_t \langle \mathbf{x}_t, \mathbf{w}_t \rangle)} \quad (2)$$

and then  $\ell(y_t | \mathbf{x}_t, \mathbf{w}_t) = -\log P(y_t | \mathbf{x}_t, \mathbf{w}_t)$ . However, in many applications (e.g., when approximating the logistic function by a Gaussian distribution) it is desirable to consider a larger class of loss function (e.g., see [25]). Since the logistic function depends only on the product  $y_t \langle \mathbf{x}_t, \mathbf{w}_t \rangle$  we set  $w_t = \langle \mathbf{x}_t, \mathbf{w}_t \rangle$  and define a function  $F(w)$  that satisfies the following properties

$$F(-w) = 1 - F(w), \quad F(0) = .5, \quad F(-\infty) = 0, \quad F(\infty) = 1. \quad (3)$$

For example,  $F$  could be the logistic function as defined in (2) or the Gaussian Cumulative Distribution Function (CDF)  $F(y_t w_t) = \Phi(\sqrt{\pi/8} y_t w_t)$  as in [1]. Then, we can re-write (2) as

$$P(y_t | \mathbf{x}_t, \mathbf{w}_t) = F(y_t \langle \mathbf{x}_t, \mathbf{w}_t \rangle), \quad (4)$$

$$\ell(y_t | \mathbf{x}_t, \mathbf{w}_t) = -\log F(y_t \langle \mathbf{x}_t, \mathbf{w}_t \rangle). \quad (5)$$

Finally, we observe that the goal of a learning algorithm is to find the best approximation  $Q(y_t | \mathbf{x}_t)$  of the unknown distribution  $P(y_t | \mathbf{x}_t, \mathbf{w}_t)$ . Therefore, we shall write  $\ell_Q(y_t | \mathbf{x}_t) = -\log Q(y_t | \mathbf{x}_t)$ . Here,  $Q$  represents an algorithm that predicts  $y_t$ .

The *pointwise regret* is defined for individual sequences  $(y_t, \mathbf{x}_t)$  as in [11], [9], [24]

$$r(y^T, Q | \mathbf{x}^T) := \sum_{t=1}^T \ell_Q(y_t | \mathbf{x}_t) - \min_{\mathbf{w}} \sum_{t=1}^T \ell(y_t | \mathbf{x}_t, \mathbf{w})$$

Thus

$$r(y^T, Q | \mathbf{x}^T) = \log \frac{\sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w})}{Q(y^T | \mathbf{x}^T)} \quad (6)$$

$$= \log \frac{\sup_{\mathbf{w}} F(y^T \langle \mathbf{x}^T, \mathbf{w} \rangle)}{Q(y^T | \mathbf{x}^T)}. \quad (7)$$

The pointwise regret  $r(y^T, Q | \mathbf{x}^T)$  is a function of  $y^T$  and  $\mathbf{x}^T$ , so it depends on *individual* sequences. Furthermore, it depends on the algorithm, represented by  $Q(y^T | \mathbf{x}^T)$ . A better measure of online logistic regression performance should decouple the regret from the fluctuations of  $y^T$  (but may still depend on the feature vector  $\mathbf{x}^T$ ) and the learning algorithm  $Q$ . Following information-theoretic view as in [4], [5], [31], we define the *maximal minimax regret* (conditioned on  $\mathbf{x}^T$ ) as

$$r_T^*(\mathbf{x}^T) := \inf_Q \max_{y^T} [r(y^T, Q | \mathbf{x}^T)]. \quad (8)$$

Notice that this definition is over all possible learning algorithms represented by  $Q$  and therefore it constitutes a lower bound for a general regret over all algorithms.

Following [26], [5] we first find a more succinct representation of the maximal minimax regret. Namely,

$$\begin{aligned} r_T^*(\mathbf{x}^T) &= \min_Q \sup_{\mathbf{w}} \max_{y^T} (-\log Q(y^T | \mathbf{x}^T) + \log P(y^T | \mathbf{x}^T, \mathbf{w})) \\ &= \min_Q \max_{y^T} [\log P^*(y^T | \mathbf{x}^T) / Q(y^T | \mathbf{x}^T)] \\ &\quad + \log \sum_{y^T} \sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w}) \\ &= \log \sum_{y^T} \sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w}) =: \log d_T(\mathbf{x}^T) \end{aligned} \quad (9)$$

if we chose  $Q(y^T | \mathbf{x}^T) = P^*(y^T | \mathbf{x}^T)$  where

$$P^*(y^T | \mathbf{x}^T) := \frac{\sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w})}{\sum_{v^T} \sup_{\mathbf{w}} P(v^T | \mathbf{x}^T, \mathbf{w})} \quad (10)$$

is the *maximum-likelihood distribution*. The above sum is often called the Shtarkov sum; see [26], [5], [10]. Observe that for not optimal  $Q$  (i.e.,  $Q \neq P^*$ ) there will be extra  $O(1)$  term in

the maximal minimax regret. We also write  $\sup_P P(y^T | \mathbf{x}^T, \mathbf{w})$  for  $\sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w})$ .

In the next section in Theorem 2 (see Theorem 4 for non-binary features) we precisely evaluate the above Shtarkov sum under additional assumption (D) (see (16) below) in which we postulate that the label probability is a function of the number (degree) of active features. But first in Theorem 1 below we show that such a regret with assumption (D) constitutes a lower bound for a general regret and arbitrary feature values.

Notice that  $r_T^*(\mathbf{x}^T)$  is still a function of the feature vector  $\mathbf{x}^T$ . To bypass this dependency, we define the *worst case maximal minimax*  $r_T^*$  as

$$r_T^* = \max_{\mathbf{x}^T} \left[ \log \sum_{y^T} \sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w}) \right], \quad (11)$$

that is,

$$r_T^* = \max_{\mathbf{x}^T} \inf_Q \max_{y^T} [r(y^T | \mathbf{x}^T)].$$

This worst case minimax regret is the closest to the minimax formulation of [20].

We may also take a probabilistic view point and assume that the feature vector is a realization of a random sequence  $\mathbf{X}^T$ . This leads to the *average maximal minimax regret* defined as

$$\bar{r}_T^* = \mathbf{E}_{\mathbf{X}^T} [r_T^*(\mathbf{X}^T)] = \mathbf{E}_{\mathbf{X}^T} [\log d_T(\mathbf{X}^T)]. \quad (12)$$

In Theorem 2 below we summarize our findings regarding the average maximal minimax regret for any learning algorithm  $Q$  and the sequence  $z^T = (y^T, \mathbf{x}^T)$ .

Finally, in the Bayesian modeling, the learning distribution is a mixture over  $\mathbf{w}$  with a prior  $\rho(\mathbf{w})$  defined as

$$Q(y^T | \mathbf{x}^T) := \int_{\mathbf{w}} \rho(\mathbf{w}) P(y^T | \mathbf{x}^T, \mathbf{w}) d\mathbf{w}.$$

In this case, the pointwise regret becomes

$$r(y^T, \rho | \mathbf{x}^T) = \log \frac{\sup_{\mathbf{w}} P(y^T | \mathbf{x}^T, \mathbf{w})}{Q(y^T | \mathbf{x}^T)}. \quad (13)$$

As in the worst case scenario, we can bypass dependency of  $\mathbf{x}^T$  by taking the maximum over  $\mathbf{x}^T$  or average over feature distribution. In Theorem 3 below we present some precise results for this Bayesian regret.

### III. MAIN RESULTS

Throughout this section, we assume that features are *binary*, that is either  $x_{i,t} = 1$  (active) or  $x_{i,t} = 0$ , unless stated otherwise. We consider two scenarios: in the *deterministic* case we assume that  $\mathbf{x}^T$  is given, while in the *stochastic* case we assume that  $\mathbf{x}^T$  is stochastically generated by some distribution. In both cases, we define  $T_j$  as the number of rounds  $t$  for which  $\sum_{i=1}^d x_{i,t} = j$ , that is, the number of feature vectors with exactly  $j$  active features. More formally,  $T_j = |\{t : \sum_{i=1}^d x_{i,t} = j\}|$ . Furthermore, by  $\alpha_j > 0$  we denote the fraction of  $\mathbf{x}_t$  ( $t = 1, \dots, T$ ) that has exactly  $j$  active features. Hence in the deterministic scenario  $\alpha_j = T_j/T > 0$  (i.e.,  $T_j > 0$ ) and  $T = T_1 + \dots + T_d$ .

In the the stochastic scenario we have

$$\alpha_j := P\left(\sum_{i=1}^d x_{i,t} = j\right) \quad (14)$$

as the probability that exactly  $j$  features are equal to 1. In particular, if we assume that  $\mathbf{x}_t = (x_{1,t}, \dots, x_{d,t})$  is distributed as the binomial( $d, p$ ) where  $P(x_{i,t} = 1) = p$ , then for all  $i \in [d]$  and all  $t$ .

$$\alpha_j := P\left(\sum_{i=1}^d x_{i,t} = j\right) = \binom{d}{j} p^j (1-p)^{d-j}. \quad (15)$$

Furthermore, in the stochastic scenario,  $(T_1, \dots, T_d)$  are random variables distributed as the multinomial  $(T, \alpha_1, \dots, \alpha_d)$ , that is,

$$P(T_1, \dots, T_d) = \binom{T}{T_1, \dots, T_d} \alpha_1^{T_1} \dots \alpha_d^{T_d}$$

where  $T = T_1 + \dots + T_d$ .

We now introduce the main assumption about  $\mathbf{w}^T$  that converts the problem from  $d$ -dimensional into one that depends only on the degree distribution (i.e., number of active features of  $\mathbf{x}^T$ ). We will assume that  $\mathbf{w}_t = (w_{1,t}, \dots, w_{d,t})$  generated according to a prior distribution are *exchangeable*. This defines our model in which the number (or degree) of occurring features determines the label probability. More precisely,

(D) For every  $k$  tuple  $(j_1, \dots, j_k) \in \{1, \dots, d\}$  and all  $1 \leq t < s \leq T$  we have

$$w^{(k)} := w_{1,t} + \dots + w_{k,t} \stackrel{d}{=} w_{j_1,s} + \dots + w_{j_k,s} \quad (16)$$

where  $\stackrel{d}{=}$  means equal ‘‘in distribution’’. In other words, the sum of  $k$  weights (those weights that have the corresponding features  $x_{i,t} = 1$ ) for any time  $t$  has the same distributions that we denote as  $w^{(k)} := w_{1,t} + \dots + w_{k,t}$ .

Observe that assumption (D) holds if all weights are *identically and independently distributed* (i.i.d.). Indeed, if say a weight distribution is  $W$ , then  $w_{1,t} + \dots + w_{k,t} \stackrel{d}{=} W \star \dots \star W$  for all  $t$  where  $\star$  denotes convolution.

Assumption (D) reformulates the problem w.r.t. the degree distribution of the feature vector. Define  $1(\mathbf{x}_t) := \sum_{i=1}^d x_{i,t}$  as the number of 1’s in  $\mathbf{x}_t$ . Then, all probabilities  $P(y_t | 1(\mathbf{x}_t) = k)$  for a given  $k$  are equal, and we denote them as

$$\theta_k(y_t) := P(y_t | 1(\mathbf{x}_t) = k) = F(y_t w^{(k)}).$$

We can view the above as a change of measure from  $\mathbf{w}$  to  $\theta$ . Finally let  $\theta_k := \theta(y_t = 1) = F(w^{(k)})$ . Then

$$P(y^T | \mathbf{x}^T, \theta) = \prod_{j=1}^d \theta_j^{k_j} (1 - \theta_j)^{T_j - k_j} \quad (17)$$

where  $T_1 + \dots + T_d = T$  and  $k_j$  is the number of  $y_t = 1$  among  $T_j$  with  $j$  active features.

Next, to estimate the maximal minimax regret, we need to compute

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

which actually becomes

$$\sup_{\theta} P(y^T | \mathbf{x}^T, \theta) = \prod_{j=1}^d \left(\frac{k_j}{T_j}\right)^{k_j} \left(\frac{T_j - k_j}{T_j}\right)^{T_j - k_j}. \quad (18)$$

This leads to the following Shtarkov sum

$$d_T(\mathbf{x}^T) = \sum_{k_1=0}^{T_1} \binom{T_1}{k_1} \left(\frac{k_1}{T_1}\right)^{k_1} \left(\frac{T_1 - k_1}{T_1}\right)^{T_1 - k_1} \dots \sum_{k_d=0}^{T_d} \binom{T_d}{k_d} \left(\frac{k_d}{T_d}\right)^{k_d} \left(\frac{T_d - k_d}{T_d}\right)^{T_d - k_d}. \quad (19)$$

This is a sophisticated sum to evaluate but we use analytic combinatorics to asymptotically evaluate it (see [28], [29]).

### A. Lower Bound

The model just introduced with assumption (D) has another interesting and useful property. It turns out that it constitutes a lower bound for the standard minimax regret  $r_T^*(\tilde{\mathbf{x}}^T)$  with general feature vector  $\tilde{\mathbf{x}}^T$ . We compare it to the minimax regret  $r_T^*(\mathbf{x}^T)$  defined above under the assumption (D) with binary features  $\mathbf{x}^T \in \{0, 1\}^T$ .

**Theorem 1.** *The minimax regret  $r^*(\mathbf{x}^T)$  with binary features under assumption (D) constitutes a lower bound for a general maximal minimax regret with features values in  $\tilde{\mathbf{x}}^T$ , that is*

$$r_T^*(\tilde{\mathbf{x}}^T) \geq r_T^*(\mathbf{x}^T). \quad (20)$$

*The regret  $r_T^*(\mathbf{x}^T)$  is precisely estimated in Theorems 2 and 4.*

*Proof.* In a general case, the probability  $P(y^T | \tilde{\mathbf{x}}^T, \mathbf{w})$  is a complicated product of probabilities that depend not only on how many active features there are but also on what these features are. Let us group all probabilities in which there is exactly one active feature under  $\theta_{1,t}(y_t | \tilde{\mathbf{x}}^t)$ , two active features under  $\theta_{2,t}(y_t | \tilde{\mathbf{x}}^t)$ , and so on until all active features under  $\theta_{d,t}(y_t | \tilde{\mathbf{x}}^t)$ . Now, to lower bound  $\sup_P P(y^T | \tilde{\mathbf{x}}^T, \mathbf{w})$  we choose particular values for  $\theta_{j,t}(y_t | \tilde{\mathbf{x}}^t)$  for each  $j$ . Namely, we set

$$\theta_{j,t}(y_t | \tilde{\mathbf{x}}^t) = \frac{k_j}{T_j}$$

where, as before,  $k_j$  is the number of  $y_t = 1$  among  $T_j$  that have  $j$  active features. But then the (log of) Shtarkov sum, and hence the maximal minimax regret becomes exactly  $r_T^*(\tilde{\mathbf{x}}^T)$  under our model (D). This completes the proof.  $\square$

### B. Precise Regret for Binary Labels

We now go back to the regret with binary features satisfying assumption (D). Using analytic combinatorics [28] and binomial sum asymptotics [14], [6], we can prove the following result regarding the asymptotic expansion of the average and worst case minimax regret under assumption (D).

**Theorem 2.** Under assumption (D) for any function  $F$  satisfying (3), the average maximal minimax regret for  $d = o(\sqrt{T})$  is given by

$$\bar{r}_T^* = \frac{d}{2} \log(T) + \frac{1}{2} \log(\alpha_1 \cdots \alpha_d) + \frac{d}{2} \log(\pi/2) + O(d/\sqrt{T}) \quad (21)$$

and its worst case minimax regret (maximum over  $\mathbf{x}^T$  or in this case  $(T_1, \dots, T_d)$ ) is

$$r_T^* = \frac{d}{2} \log\left(\frac{T}{d}\right) + \frac{d}{2} \log(\pi/2) + O(d/\sqrt{T}) \quad (22)$$

for large  $T$ .

Let us now present results for the Bayesian pointwise regret as defined in (13) which becomes

$$r(y^T, \rho | \mathbf{x}^T) = \log \frac{\sup_{\theta_1, \dots, \theta_d} \prod_{j=1}^d \theta_j^{k_j} (1 - \theta_j)^{T_j - k_j}}{\int_{\theta} \rho(\theta) \prod_{j=1}^d \theta_j^{k_j} (1 - \theta_j)^{T_j - k_j}}. \quad (23)$$

By (18) we can re-write it as (with  $T_i > 0$ )

$$r(y^T, \rho | \mathbf{x}^T) = \log \frac{\prod_{j=1}^d \binom{k_j}{T_j} \binom{T_j - k_j}{T_j}^{T_j - k_j}}{\int_{\theta} \rho(\theta) \prod_{j=1}^d \theta_j^{k_j} (1 - \theta_j)^{T_j - k_j}}. \quad (24)$$

Using Stirling's approximation, Dirichlet distribution, binomial sum asymptotics [3], [6], [14] and analytic combinatorics [28] we prove the following results.

**Theorem 3.** For the Bayesian setting with Jeffrey's prior presented under assumption (D) for any function  $F$  satisfying (3), the pointwise regret in the deterministic setting becomes

$$r(y^T | \mathbf{x}^T) = \frac{1}{2} \log(T_1 \cdots T_d) + \frac{d}{2} \log(\pi/2) + O(1/T) \quad (25)$$

for large  $T_j$  and  $k_j$ . The average pointwise regret is

$$\mathbf{E}_{\mathbf{x}}[r(y^T | \mathbf{x}^T)] = \frac{d}{2} \log T + \frac{1}{2} \log(\alpha_1 \cdots \alpha_d) + \frac{d}{2} \log(\pi/2) + O(d/T) \quad (26)$$

for large  $T$  with  $\alpha_i > 0$ .

### C. Regret for Non-binary Labels

We briefly discuss how to extend our results to non-binary labels, say, label alphabet  $\mathcal{Y}$  of size  $m$ . Following [9] we need to extend the weight vector to the weight matrix  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{m-1}]$  such that  $\mathbf{w}_i = (w_{1,i}, \dots, w_{d,i})$ . Then the multinomial logistic function known also as *softmax function* is defined as

$$P(y_t = \ell | \mathbf{x}_t, \mathbf{W}_t) = \frac{e^{\langle \mathbf{x}_t, \mathbf{w}_{t,\ell} \rangle}}{\sum_{k=1}^m e^{\langle \mathbf{x}_t, \mathbf{w}_{t,k} \rangle}} \quad (27)$$

for  $\ell = 0, 1, \dots, m-1$ . Also, as before we define

$$\theta_{k,\ell}(y_t) := P(y_t = \ell | \mathbf{x}_t = k, \mathbf{W})$$

leading to

$$P(y^T | \mathbf{x}^T, \theta) = \prod_{j=1}^d \prod_{\ell=1}^m \theta_{j,\ell}^{k_{j,\ell}} \quad (28)$$

where  $k_{j,\ell}$  is the number labels equal to  $\ell$  among  $T_j$ .

Following the footsteps of our analysis for the binary labels we arrive at the following generalization of Theorem 2.

**Theorem 4.** Under assumption (D) for any function  $F$  satisfying (3) the average maximal minimax regret for  $m = O(1)$  and  $d = o(\sqrt{T})$  becomes

$$\bar{r}_T^* = \frac{d(m-1)}{2} \log(T/2) + \frac{(m-1)d}{2} \log(\alpha_1 \cdots \alpha_d) + \frac{d}{2} \log(\pi/\Gamma^2(m/2)) + O(d/\sqrt{T}) \quad (29)$$

(with  $\Gamma(x)$  being the Euler gamma function) while for  $dm^{3/2} = o(\sqrt{T})$  we find

$$\bar{r}_T^* = \frac{d(m-1)}{2} \log(T/m) + \frac{(m-1)d}{2} \log(\alpha_1 \cdots \alpha_d) + \frac{md}{2} \log e + O(dm^{3/2}/\sqrt{T}). \quad (30)$$

Furthermore, the worst case maximal minimax for  $dm^{3/2} = o(\sqrt{T})$  is

$$r_T^* = \frac{d(m-1)}{2} \log(T/(dm)) + \frac{md}{2} \log e + O(dm^{3/2}/\sqrt{T}) \quad (31)$$

for large  $T$ .

Finally, we go back to the lower bound discussed in Theorem 1. We compare the pointwise regret  $r_T(y^T, Q | \tilde{\mathbf{x}}^T)$ ; with general feature  $\tilde{\mathbf{x}}^T$ , to that in Theorem 4; the worst case minimax regret  $r_T^*$  (maximized over  $\mathbf{x}^T \in \{0, 1\}^T$ ). We conclude the following general lower bound.

**Corollary 1.** Consider a general pointwise regret  $r_T(y^T, Q | \tilde{\mathbf{x}}^T)$  for any algorithm/ distribution  $Q$  and general feature values. The following holds for  $dm^{3/2} = o(\sqrt{T})$

$$\max_{(\tilde{\mathbf{x}}^T, y^T)} r_T(y^T, Q | \tilde{\mathbf{x}}^T) \geq \frac{d(m-1)}{2} \log(T/(dm)) + \frac{md}{2} \log e + O(dm^{3/2}/\sqrt{T}) \quad (32)$$

for large  $T$ .

## IV. CONCLUSION

We described a logistic regression problem in which the label probability is determined by the degree of existence of features in an example, instead of by the actual features. Online regret for this problem, while interesting by itself for applications such as graph models, and specifically for social networks, can also be used to bridge well established results from the universal compression literature to the study of regret in machine learning. We demonstrated that for this novel problem we can precisely compute various variants of the regret, showing logarithmic regret, which linearly increases with the dimensionality. The precise regret for this problem by itself serves as a *general lower bound* for the regret of standard logistic regression. This connection opens up a large range of possibilities to apply established theory in the study of universal compression redundancy to studying regret for online learning problems.

## ACKNOWLEDGMENT

This work was done while Wojciech Szpankowski was visiting Google.

## REFERENCES

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, New York, 2006.
- [2] Y. Choi and W. Szpankowski. Compression of graphical structures: Fundamental limits, algorithms, and experiments. *IEEE Trans. Information Theory*, 58:620–638, 2012.
- [3] J. Cichon and Z. Golebiewski. On the bernoulli sums and bernstein polynomials. In *AofA*, pages 179–190, 2012.
- [4] L. D. Davisson. Universal noiseless coding. *IEEE Trans. Inf. Theory*, IT-19(6):783–795, Nov. 1973.
- [5] M. Drmota and W. Szpankowski. Precise minimax redundancy and regrets. *IEEE Trans. Inf. Theory*, IT-50:2686–2707, 2004.
- [6] P. Flajolet. Singularity analysis and asymptotics of bernoulli sums. *Theoretical Computer Science*, 215:371–381, 1999.
- [7] P. Flajolet and R. Sedgewick. *Analytic Combinatorics*. Cambridge University Press, 2008.
- [8] P. Flajolet and W. Szpankowski. Analytic variations on redundancy rates of renewal processes. *IEEE Trans. Information Theory*, 48:2911–2921, 2002.
- [9] Dylan J Foster, Satyen Kale, Haipeng Luo and Mehryar Mohri, and Karthik Sridharan. Logistic regression: The importance of being improper. In *COLT - Conference on Learning Theory*, 2018.
- [10] P. D. Grunwald. *The Minimum Description Length Principle*. MIT Press, 2007.
- [11] E. Hazan. The convex optimization approach to regret minimization. In S. Sra, S. Nowozin, and S. Wright, editors, *Optimization for Machine Learning*, pages 287–303. MIT press, 2012.
- [12] E. Hazan, T. Koren, and K. Y. Levy. Logistic regression: Tight bounds for stochastic and online optimization. In *The 27th Conference on Learning Theory, COLT 2014*, page 197–209. MIT press, 2014.
- [13] P. Jacquet, G. I. Shamir, and W. Szpankowski. Precise minimax regret for logistic regression with categorical feature values, 2020.
- [14] P. Jacquet and W. Szpankowski. Entropy computations via analytic depoissonization. *IEEE Trans. on Information Theory*, 45:1072–1081, 1999.
- [15] R. Jezequel, P. Gillard, and A. Rudi. Efficient improper learning for online logistic regression. In *PMLR, COLT’20*, volume 125, 2020.
- [16] Sham M Kakade and Andrew Y. Ng. Online bounds for bayesian algorithms. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 641–648. MIT Press, 2005.
- [17] Tomasz Łuczak, Abram Magner, and Wojciech Szpankowski. Asymmetry and structural information in preferential attachment graphs. *Random Structures & Algorithms*, 55(3):696–718, 2019.
- [18] H. B. McMahan and M. J. Streeter. Open problem: Better bounds for online logistic regression. In *Journal of Machine Learning Research-Proceedings Track*, 23, 2012.
- [19] A. Orlitsky and N. P. Santhanam. Speaking of infinity. *IEEE Trans. Inf. Theory*, 50(10):2215–2230, Oct. 2004.
- [20] A. Rakhlin and K. Sridharan. Online nonparametric regression with general los function. In *COLT*, 2014.
- [21] J. Rissanen. Fisher information and stochastic complexity. *IEEE Trans. Information Theory*, 42:40–47, 1996.
- [22] S. Shalev-Schwartz and S. Ben-David. *Understanding Machine learning*. Cambridge University Press, 2014.
- [23] G. I. Shamir. On the MDL principle for i.i.d. sources with large alphabets. *IEEE Trans. Inform. Theory*, 52(5):1939–1955, May 2006.
- [24] G. I. Shamir. Logistic regression regret: What’s the catch? In *COLT*, 2020.
- [25] G. I. Shamir and W. Szpankowski. Low complexity approximate bayesian logistic regression for sparse online learning, 2020.
- [26] Y. M. Shtarkov. Universal sequential coding of single messages. *Problems of Information Transmission*, 23(3):3–17, Jul.-Sep. 1987.
- [27] W. Szpankowski. On asymptotics of certain recurrences arising in universal coding. *Problems of Information Transmission*, 34:55–61, 1998.
- [28] W. Szpankowski. *Average Case Analysis of Algorithms on Sequences*. Wiley, New York, 2001.
- [29] W. Szpankowski and M. Weinberger. Minimax pointwise redundancy for memoryless models over large alphabets. *IEEE Trans. Information Theory*, 58:4094–4104, 2012.
- [30] Q. Xie and A. Barron. Minimax redundancy for the class of memoryless sources. *IEEE Trans. Information Theory*, pages 647–657, 1997.
- [31] Q. Xie and A. Barron. Asymptotic minimax regret for data compression, gambling, and prediction. *IEEE Trans. Information Theory*, 46:431–445, 2000.