

# A Universal Probability Assignment for Prediction of Individual Sequences

Yuval Lomnitz, Meir Feder

Tel Aviv University, Dept. of EE-Systems

Email: yuval.lomnitz@gmail.com ,meir@eng.tau.ac.il

**Abstract**—Is it a good idea to use the frequency of events in the past, as a guide to their frequency in the future (as we all do anyway)? In this paper the question is attacked from the perspective of universal prediction of individual sequences. It is shown that there is a universal sequential probability assignment, such that for a large class loss functions (optimization goals), the predictor minimizing the expected loss under this probability, is a good universal predictor. The proposed probability assignment is based on randomly dithering the empirical frequencies of states in the past, and it is easy to show that randomization is essential. This yields a very simple universal prediction scheme which is similar to Follow-the-Perturbed-Leader (FPL) and works for a large class of loss functions, as well as a partial justification for using probabilistic assumptions.

## I. INTRODUCTION

In this paper the problem of universal sequential prediction of an individual unknown sequence is considered [1][2][3], and a prediction approach based on universal probability assignment is proposed. Given a space of strategies  $\mathcal{B}$ , a space of nature states  $\mathcal{X}$  and a loss function  $l(b, x), b \in \mathcal{B}, x \in \mathcal{X}$ , the purpose is to assign the next strategy  $\hat{b}_t$  given the knowledge of the past states  $\mathbf{x}_1^{t-1}$ , such that the overall loss  $\sum_{t=1}^n l(\hat{b}_t, x_t)$  would be asymptotically close to the loss obtained by the best fixed strategy known a-posteriori after viewing the entire sequence  $\mathbf{x}_1^n$ , i.e.  $\min_{b \in \mathcal{B}} \sum_{t=1}^n l(b, x_t)$ .

In the particular case of sequential probability assignment under the log loss function  $l(b, x) = \log \frac{1}{b(x)}$  where  $\mathcal{B}$  is the space of probability assignments on the finite alphabet  $\mathcal{X}$ , or equivalently in universal sequential compression, it is shown [3][4, §13][1, §9] that it is possible to assign probabilities  $\hat{p}_t(x_t)$  for the next state in an arbitrary sequence of states  $x_t \in \mathcal{X}, t = 1, 2, \dots, n$ , given the past states, such that for any possible sequence, the overall probability  $\hat{p}(\mathbf{x}) = \prod_{t=1}^n \hat{p}_t(x_t)$  would not be too far, in a multiplicative or logarithmic sense, from the best i.i.d. probability assigned to the sequence a-posteriori  $\max_{p(\cdot)} \prod_{t=1}^n p(x_t)$ . The result extends to probability assigned by Markov machines or finite state machines [5]. This problem is related to universal compression because the overall compression length corresponds to  $\log \left( \frac{1}{p(\mathbf{x})} \right)$ . A remarkable feature of these universal probability assignments is that, although nothing is assumed about the sequence, to construct a universal encoder it is enough to encode as if  $\hat{p}_t(\cdot)$  was the *true* probability of the next state.

These universal probability assignments, such as the Laplace [4, §13.2] or Krichevsky-Trofimov (KT) [6] estimators, have an intuitively appealing structure which induces

a small bias over the empirical distribution seen so far. For example, Laplace's estimate for the probability distribution of  $x_t$  is

$$\hat{p}_t(x) = \frac{N_{t-1}(x) + 1}{(t-1) + |\mathcal{X}|}, \quad (1)$$

where  $N_t(x)$  denotes the number of times the state  $x$  appears in  $\mathbf{x}_1^t$ . While these estimators get closer with time to the measured empirical distribution, they do not “trust” it completely, and, for example, never assign a probability value 0 to states that had not appeared before. Furthermore, in the probabilistic prediction setting the same distributions were shown to perform well not only for the log loss: the predictor which minimizes the expected loss under these distributions  $\hat{b}_t = \operatorname{argmin}_b \mathbb{E}_{X \sim \hat{p}_t(\cdot)} l(b, X)$  operates well for a wider class of loss functions [3, §III.A.2].

This naturally leads to the following question: is it possible to forecast an individual sequence by first generating a probability assignment based on the past, and then minimizing the expected loss under this assignment (i.e. in a way, acting as if future events truly happen with this probability)? Consider prediction schemes of the following form:

- 1) Generate a probability assignment  $P_t^{(u)}(x)$  based on the past of the sequence  $\mathbf{x}_1^{t-1}$ , in a way which does not depend on the loss function.
- 2) To predict  $b_t$  under the loss function  $l(b, x)$ , choose the strategy that minimizes the expected loss under  $P_t^{(u)}$ , i.e.:

$$\hat{b}_t = \operatorname{argmin}_b \mathbb{E}_{X \sim P_t^{(u)}(\cdot)} [l(b, X)] \quad (2)$$

If there exists a single scheme for generating  $P_t^{(u)}(\cdot)$  that does not depend on the loss function  $l(b, x)$ , but for which  $\hat{b}_t$  yields a good (Hannan-consistent [1]) predictor for a certain class of loss functions, then we call  $P_t^{(u)}(\cdot)$  a *universal sequential probability assignment* with regards to that class. Notice that this term has been used in the past with respect to the log-loss, so the definition above can be considered a natural extension.

It is easy to show that, if the class of loss functions includes even simple loss functions such as the 0-1 loss (the number of errors), then no deterministic assignment can be universal, and therefore the Laplace or KT assignments are inadequate. However, it is shown in this paper that the random assignment obtained by slightly perturbing the empirical frequencies is universal for a large class of loss functions, including the log-loss and any bounded loss.

In addition to supplying a simple and general universal prediction scheme, this result also has interpretations contributing to our understanding of probability. For example, it supplies justification for treating the statistics of a process in the past as a guide to its statistics in the future, without having to assume the process is indeed stationary, or that it is driven by a “probabilistic” law. In other words, if our natural behavior is in some way similar to the prediction algorithm described here, then the claims on its convergence can be used to justify this behavior.

The next section completes the problem definition and discusses the boundaries of the solution, and relations to known results. Section III gives the main results, and Section IV discusses the possible implications on understanding probabilistic behavior. The proofs are given in Section V.

## II. PROBLEM STATEMENT AND DISCUSSION

Building upon the definitions already presented in the introduction, in this section some complementary definitions are presented. We assume throughout this paper that  $\mathcal{X}$  is finite (otherwise there is no meaning to measuring empirical frequencies). The set of possible strategies  $\mathcal{B}$  is not restricted. The loss function  $l(b, x)$  is constant over time.

Let us define the accumulated loss of a sequential predictor  $\hat{b}_t(\mathbf{x}_1^{t-1})$  as:

$$\hat{L}_n = \sum_{t=1}^n l(\hat{b}_t, x_t), \quad (3)$$

and the loss of the best fixed strategy as:

$$L_n^* = \min_b \sum_{t=1}^n l(b, x_t). \quad (4)$$

The difference  $\hat{L}_n - L_n^*$  which is defined as the regret, is a function of the predictor and the sequence. The worst case regret is:

$$\mathcal{R}_{\max} = \max_{\mathbf{x}_1^n} (\hat{L}_n - L_n^*), \quad (5)$$

and the normalized regret is  $\frac{\mathcal{R}_{\max}}{n}$ . A forecasting strategy  $\hat{b}_t$  is said to be *Hannan-consistent*, if  $\limsup_{n \rightarrow \infty} \frac{\mathcal{R}_{\max}}{n} \leq 0$  almost surely (the probability is over the randomization in the forecaster if it is random). This means that for large  $n$ , the loss of the forecaster is essentially at least as small as that of any fixed strategy. As mentioned in the introduction, the problem addressed in this paper is of finding a sequential probability assignment  $P_t^{(u)}(\cdot)$  such that the resulting prediction scheme (2) is Hannan-consistent for a large class of loss functions. We will focus mainly on bounding the *expected* loss (over the predictor’s randomization), because it also leads to almost-sure bounds by applying the strong law of large numbers. The maximum *expected* regret is defined as:

$$\overline{\mathcal{R}}_{\max} = \max_{\mathbf{x}_1^n} \mathbb{E} [\hat{L}_n - L_n^*], \quad (6)$$

For some loss functions satisfying smoothness conditions [1, Thm 3.1][2, Thm 1], the forecasting strategy known as “Follow the Leader” (FL), which chooses at each time the best strategy in retrospect  $\hat{b}_t^{(\text{FL})} = \operatorname{argmin}_b \sum_{i=1}^{t-1} l(b, x_i)$ ,

is Hannan consistent. Rewriting the above as  $\hat{b}_t^{(\text{FL})} = \operatorname{argmin}_b \sum_{x \in \mathcal{X}} \frac{N_{t-1}(x)}{t-1} l(b, x)$ , it can be interpreted as an implementation of (2) where the universal probability assignment equals the empirical frequencies  $P_t^{(u)}(\cdot) = \frac{N_{t-1}(x)}{t-1}$ . In other words, for this family of loss functions, there is a simple solution for  $P_t^{(u)}(\cdot)$ , namely the empirical distribution. However this class of loss functions where FL is universal, is rather limited.

For a probability assignment to be “general” enough, one would want to cover, at the least, the family of discrete-strategy, discrete-state loss functions, presented by Hannan [7]. For this family, the loss function can be represented by a general  $|\mathcal{B}| \times |\mathcal{X}|$  matrix specifying the loss for each strategy and each state of nature. It is well known [1, §4] and straightforward to see that randomization is required in order to cover this class: consider the 0-1 loss case, i.e. binary sequences  $\mathcal{X} = \mathcal{B} = \{0, 1\}$  with  $l(b, x) = \text{Ind}(b \neq x)$ , where the total loss is the number of errors. For this loss function, no deterministic predictor yields Hannan-consistency, because for each deterministic predictor there exists a sequence which fails the predictor completely, by choosing the next outcome as the opposite of the predictor’s choice, while the loss of the best fixed predictor is at most  $n/2$ . Because a deterministic  $P_t^{(u)}(\cdot)$  inevitably leads to a deterministic predictor (2), this implies a random  $P_t^{(u)}(\cdot)$  is required, in general.

For the binary 0-1 loss problem, Feder, Merhav and Gutman [8] used a small dither when the empirical probability is close to  $\frac{1}{2}$ , which effectively avoids a decision when the frequencies of 0, 1 are nearly equal.<sup>1</sup> For this specific problem, the optimal solution (in the sense of minimax regret) is known exactly and was presented by Cover [9]. While the optimal dither in this problem is different than the straight line used by Feder, Merhav and Gutman, and is not known in general, this is of no consequence in the current problem, as we are only considering Hannan consistency. This solution, as well as the small bias from the empirical distribution which is required in the log-loss problem (1), motivates the following choice of  $P_t^{(u)}(\cdot)$ : add a small dither to  $N_{t-1}(x)$  (the counts of events in the past) and re-normalize. As shown below, this solution achieves Hannan-consistency for any bounded loss function and for the log loss.

The proposed forecaster is reminiscent of the scheme termed “Follow the Perturbed Leader” (FPL), originally proposed by Hannan [7], in which the decision is obtained by adding a small dither to the accumulated loss of every reference strategy and then choosing the best one. Indeed, dithering the frequencies is similar, but not equivalent, to dithering the accumulated losses, and our proof technique for the bounded loss case borrows from Kalai and Vempala’s [10]. Following this similarity we term the scheme proposed here “Follow the Perturbed Frequency” (FPF). Notice, however, that FPL is defined, in general, only when the number of strategies is finite, while FPF is defined, in general, only when the number of outcomes (states) is finite, and does not have to assume

<sup>1</sup>It is interesting to note that for the 0-1 loss problem their forecaster is equivalent to a “Follow the Perturbed Leader” forecaster with a uniform distribution (see below) and also equivalent to the forecaster proposed here.

the number of strategies is finite. On the other hand, FPL can deal with more general forms of the problem, including time-varying loss functions.

The problem considered here is a close relative of the calibration problem [1, §4.5], i.e. the problem of estimating from an individual sequence, probability forecasts that pass certain consistency tests. The problems are related in that, in both cases it is shown possible to generate from empirical data collected from an individual sequence, probability assignments that appear to operate as well as forecasts which are based on knowledge of the “true” statistical model. Also, randomization is essential in both cases. However, none of the problems is a special case of the other: the probability assignment shown here is not necessarily calibrated, and a calibrated probability assignment does not necessarily satisfy the requirements of the current problem.<sup>2</sup>

In this paper, in order to simplify matters, only fixed strategies are considered. As one of our motivations is to rationalize the behavior of learning probabilities from the past, it is enough to consider fixed strategies in order to see the advantage of this behavior. The extension to dynamic reference strategies is unfortunately not immediate as in the setting of prediction with expert advice [1, §2], where dynamic strategies can be turned into fixed ones by simple enumeration (i.e. replacing the strategy with the index of the strategy), because we explicitly assume a fixed loss function. However in some cases, the core of the prediction problem lies in competing with fixed strategies. For example, reference strategies defined by states (such as Markov predictors or finite state machines), can be considered as fixed strategies in each sub-sequence belonging to the same state.

### III. MAIN RESULTS

Let  $N_t(x)$  be number of times a specific  $x$  occurred in the sequence  $\mathbf{x}$  up to and including time  $t$ . The universal sequential probability assignment is defined as:

$$\begin{aligned} P_t^{(u)}(x) &= c_t \cdot (N_{t-1}(x) + h_t \cdot u_t(x)) \\ &= \frac{N_{t-1}(x) + h_t \cdot u_t(x)}{t - 1 + h_t \cdot \sum_{x' \in \mathcal{X}} u_t(x')} \end{aligned} \quad (7)$$

where  $c_t = \sum_{x \in \mathcal{X}} (N_{t-1}(x) + h_t \cdot u_t(x))$  is the normalizer guaranteeing unit sum.  $u_t(x) \sim U[0, 1]$  is a random dither which is assumed to be uniformly distributed, i.i.d. over different  $x$  and  $t$  (dependence over  $t$  does not affect the expected regret).  $h_t$  is a non-decreasing positive sequence. Our philosophical considerations (i.e. justifying probabilistic behavior) motivate keeping  $h_t$  as general as possible rather than finding a specific optimal sequence  $h_t$  for each problem.

The FPF predictor, for any loss function  $l(b, x)$  is defined by:

$$b_t^{(\text{FPF})} = \operatorname{argmin}_{b \in \mathcal{B}} \mathbb{E}_{\mathbf{X} \sim P_t^{(u)}(x)} [l(b, X)] \quad (8)$$

<sup>2</sup>Consider for example the 0-1 loss problem, and a sequence containing an equal number of zeros and ones. Any probability forecaster yielding only values in the range  $0.5 \pm \epsilon$  is  $\epsilon$ -calibrated, while the decisions based on these probabilities (when plugged into (2)) can be arbitrary (depending on whether the probability is smaller or larger than 0.5), and can yield arbitrarily bad (or good) aggregate losses.

**Theorem 1.** Assuming  $h_t = h_1 \cdot t^\alpha$ , with  $\alpha \in (0, 1)$ , the FPF predictor is Hannan-consistent for any bounded loss function and for the log-loss. Therefore under these conditions,  $P_t^{(u)}(x)$  defined in (7) is a universal probability assignment for the class.

This theorem is based on the two following theorems:

**Theorem 2.** Assume the loss function is bounded  $|l(b, x)| \leq R$ . Then:

- 1) The expected regret of FPF is upper bounded by

$$\overline{\mathcal{R}}_{\max} \leq 2R \sum_{t=1}^n h_t^{-1} + 2R|\mathcal{X}|h_n \quad (9)$$

- 2) Particularly, for any  $h_t = h_1 \cdot t^\alpha$ , with  $\alpha \in (0, 1)$ , the normalized expected regret  $\frac{1}{n} \overline{\mathcal{R}}_{\max}$  tends to zero with  $n$ .
- 3) For  $h_t = \sqrt{\frac{2t}{|\mathcal{X}|}}$ ,  $\frac{1}{n} \overline{\mathcal{R}}_{\max} \leq 4R\sqrt{\frac{2|\mathcal{X}|}{n}}$ .

**Corollary 2.1.** The theorem holds under a milder condition, that the loss function is bounded only for the set of optimizing strategies, defined as

$$\mathcal{B}_{\text{opt}} = \left\{ \operatorname{argmin}_{b \in \mathcal{B}} \sum_{x \in \mathcal{X}} \lambda(x) l(b, x) : \lambda(x) \geq 0, \exists x : \lambda(x) > 0 \right\} \quad (10)$$

and where  $R = \sup_{x \in \mathcal{X}, b \in \mathcal{B}_{\text{opt}}} l(b, x)$ . Particularly, the theorem holds for the  $L_2$  norm loss,  $l(\mathbf{b}, \mathbf{x}) = \|\mathbf{b} - \mathbf{x}\|^2$  for  $\mathcal{X} \subset \mathbb{R}^d$  ( $|\mathcal{X}| < \infty$ ), and  $\mathcal{B} = \mathbb{R}^d$ . In that case  $R = \max_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} \|\mathbf{x} - \mathbf{x}'\|^2$  is the squared diameter of the set  $\mathcal{X}$ .

Notice that in the most general case without any limitations (such as on magnitude), it is generally impossible to devise a universal scheme for the  $L_2$  norm loss that beats the best fixed strategy, i.e. the empirical mean up to a constant, and it is made possible in the current problem by the assumption that  $\mathcal{X}$  is finite.

The proof of Theorem 2 is similar in spirit to the proof of Kalai and Vempala [10] for the FPL forecaster, as the perturbation on  $N_t(x)$  can be translated to a perturbation on the accumulated loss.

**Theorem 3.** For the case of the log-loss, where  $b(x)$  is a probability distribution over  $\mathcal{X}$  and  $l(b, x) = \log\left(\frac{1}{b(x)}\right)$ ,<sup>3</sup> the expected regret of FPF satisfies:

- 1)

$$\overline{\mathcal{R}}_{\max} \leq \sum_{t=1}^n \frac{|\mathcal{X}|h_t - 1}{t} + \sum_{t=1}^n \frac{1}{\lfloor \frac{t-1}{|\mathcal{X}|} \rfloor + h_t} \quad (11)$$

- 2) Particularly, for any  $h_t = h_1 \cdot t^\alpha$ , with  $\alpha \in [0, 1)$ , the normalized expected regret  $\frac{1}{n} \overline{\mathcal{R}}_{\max}$  tends to zero with  $n$ .
- 3) For constant  $h_t$  the expected regret behaves like  $O(\log n)$  and specifically for the choice  $h = |\mathcal{X}|^{-1}$ ,  $\overline{\mathcal{R}}_{\max} \leq |\mathcal{X}| \log(n)$ .

<sup>3</sup>All log-s in this paper are in the natural base.



Regarding the last case, notice that this redundancy is similar to the redundancy obtained with Laplace's estimator and approximately twice the redundancy obtained using Krichevsky-Trofimov's (which is approximately  $\frac{|\mathcal{X}|-1}{2} \log n$ ). However notice that the target of the FPF forecaster was not to produce optimal redundancy for specific loss functions.

The proofs of the theorems stated above appear in Section V below.

#### IV. IMPLICATIONS ON THE UNDERSTANDING OF PROBABILITY

##### A. Initial probabilities

A basic question in the application and philosophy of probability theory is: where do initial probabilities originate from (see, e.g. [11]) ? The fact is, that in many situations a probability distribution is deduced from the relative frequency of events in the past. While this deduction may be justified based on some stationarity assumption, it is often used exactly in those situations where precise analysis of the source of events is not possible, and therefore the assumption that the frequency of events in the future would be similar to their frequency in the past is not necessarily justified. In spite of this, we often deduce a probability distribution based on past statistics and use this probability for decision making with regards to future events. It seems that not only humans but also animals use this principle [12].

One motivation for the problem posed in Section II, of searching for a universal probability assignment, is the attempt to justify this behavior based on mathematical, rather than physical assumptions. The theory of universal prediction of individual sequences, or repetitive games, seems a good framework for this purpose, because it facilitates deduction from the past, without assumptions that the past indicates anything with respect to the future. The existing universal prediction schemes are less suitable for this purpose since they determine the next strategy in a contrived way, as a function of the past frequencies and the loss function, whereas in the probability-based decision making, it is assumed that there exist a single "true" probability.

The success of the FPF predictor for a large set of loss functions, indicates that indeed it is useful to rely on past frequencies, and draw from them a "probability" distribution, even if the future is arbitrary. The dither may be interpreted as the assumption that the future would be similar but not identical to the past, and prevents using a too "decisive" strategy (such as choosing '0' or '1' in the 0-1 loss case), based on a small change in the frequencies. It would be farfetched to claim that this is *the* justification for using probabilities: clearly the reason is related to the regularity that many natural processes exhibit; however it supplements our intuitive understanding by showing that even if these assumptions fail, there is still benefit in learning probabilities from the past.

##### B. Meaning of probability

In the previous section we tried to justify a specific choice of a probability. However, probability itself is not a well

defined concept, and many attempts to explain or justify its use have been made. A good introduction to these philosophical questions can be found in [11] (for a quick overview see [13, Chap. ??]). While there is no dispute on the mathematical axiomatic theory dealing with probability functions, the meaning of probability, and the justification for using it are questionable.

In a nutshell, the main interpretations to probability are the relative frequency approach, a-priori or logical approach and the subjectivistic approach. Relative frequency theories interpret probability as the limiting frequencies in very large groups of events (called "collectives"). A-priori theories interpret probability as logical relation between sentences, and an extension of formal logic: the attributes "true" and "false" are represented by probabilities of 1 and 0, and are extended by adding a range of probabilities in between. Subjectivistic theories interpret probability as a measure of the degree of belief of a certain person in a certain proposition, and therefore its value is not unique.

A main issue in all interpretations is what probability means with respect to the future. The current results can be interpreted under the framework of the subjectivistic theories, which view probability as a tool for decision making, i.e. probability is just the relative weight that we put on each future event when making decisions. Because under subjectivistic theories any probability is valid, there is a problem of justifying any specific choice of a probability assignment, as well as the merit of making decisions according to probabilistic considerations.

The current results can be thought of as a partial resolution to this question: the suggestion of learning probability from the past by biasing or dithering past frequencies, is a good one in the sense that it is better than any fixed behavior (and as a result, of making decisions according to any fixed probability). This demonstrates a clear merit in following probabilistic considerations, which is not dependent on any assumptions with respect to the real world (the process  $\mathbf{x}_t$ ).

There are some issues, however, with this interpretation. First, the problem setting is limited, compared to our actual use of probability. Learning from experience extends far beyond the framework of repetitive games and constant loss functions, as we usually deduce probabilities from the past and use them to solve *new* problems. Also the fact  $\mathcal{X}$  is assumed discrete is somewhat limiting, although it may be sufficient to justify probabilistic intuition, which is fundamentally based on distributions on finite sets (such as coins and dice).

But the main weakness of this interpretation is that it relies on randomness for generating the universal probability assignment  $P^{(u)}$  (and as a result, the claims we can make are also probabilistic), and so it may lead to a cyclic argument of explaining probability by using probability. The randomness used here is in a restricted form of "controlled randomness" which is generated by the forecaster. I.e. if we believe it is possible to draw random coins, it is enough for this interpretation to hold and be meaningful. An alternative assumption is pseudo-randomness, i.e. assume that we can generate the dither not randomly, but such that "nature" (drawing the next  $x_t$ ) cannot guess it, and it appears effectively random.

Unfortunately, like in many other theories, we are not able to escape some form of “belief” or conjecture with respect to the future.

Another way to avoid the need for randomness is to avoid problems such as the 0-1 loss case, in which one is forced to bet, problems that are insolvable without randomness. For example, if the loss is convex with respect to the strategy, then the loss when taking the expected value of a random strategy  $b$  is always better than the expected loss when  $b$  is random. In this case, the forecaster can make a deterministic decision: replace (8) with  $\hat{b}_t = \mathbb{E} \left\{ \hat{b}_t^{(\text{FPF})} \right\}$ , where the expected value is with respect to the randomness of  $P^{(u)}$ . This can be thought of as a different rule for making decisions based on the past: take as probability the empirical frequencies in the past, however when making a decision which changes significantly with respect to small variations in the probability, take the average decision over these small variations. This rule is deterministic and aligns with intuition, however the restriction to “smooth” loss functions may be too limiting.

Another question that would naturally arise with respect to this explanation is how it aligns with the fact that, at least in the theoretical application of probability theory (e.g. estimation theory, communication theory) we do not use dithers in our probabilities. It seems that the idea of dithering the probabilities is a similar notion to the idea of checking sensitivity of a given solution to the probabilistic assumptions. In case the solution to a given problem does not depend crucially on the exact probability values, adding the dither is indeed redundant. On the other hand, if the solution depends crucially on a small change in the probabilistic assumptions, it may be reasonable to doubt its operation in the real world.

## V. PROOFS

### A. Proof of Theorem 2

The proof follows the same line of thought of Kalai-Vempala [10]: first, the regret of a clairvoyant forecaster using  $x_t$  in addition to  $\mathbf{x}_1^{t-1}$  is bounded. Then, the difference in performance between the clairvoyant forecaster and the proposed forecaster is bounded, by using the fact that some of the dither works in the same direction as the the difference between them.

1) *Definitions:* The cumulative loss for playing the constant strategy  $b$  up to time  $t$  is  $L_t(b) = \sum_{i=1}^t l(b, x_i)$ . We denote for brevity  $\mathbf{B}\{L(b)\} \triangleq \argmin_b L(b)$  the best strategy for cumulative loss function  $L(b)$ .

The optimal fixed (a-posteriori) best fixed strategy is  $\mathbf{B}\{L_n(b)\}$  and has loss  $L_n^* = L_n(\mathbf{B}\{L_n(b)\})$ . As another example to clarify the notation, the FL predictor can be written as  $\mathbf{B}\{L_{t-1}(b)\}$ , and the FPL predictor [10] can be written  $\mathbf{B}\{L_{t-1}(b) + p_t(b)\}$  where  $p_t(b)$  is a random perturbation.

Let us define the dithered count at time  $t-1$  as

$$N_{t-1}^{(p)}(x) \triangleq N_{t-1}(x) + h_t u_t(x), \quad (12)$$

and the respective dithered accumulated loss as

$$L_{t-1}^{(p)}(b) \triangleq \sum_x l(b, x) N_{t-1}^{(p)}(x). \quad (13)$$

This loss could be thought of as the loss during a sequence which is an extension of the actual sequence with some random states. Notice the distinction between  $\hat{L}_n$  defined in (3), which is the loss of the universal predictor, and  $L_t^{(p)}$  which is the accumulated loss whose minimization yields the predictor. The distribution  $P^{(u)}$  is proportional to  $N_{t-1}^{(p)}(x)$ , and thus the FPF forecaster is equivalent to optimizing the dithered loss:

$$\begin{aligned} b_t^{(\text{FPF})} &= \argmin_{b \in \mathcal{B}} \mathbb{E}_{\mathbf{x} \sim P_t^{(u)}(x)} [l(b, X)] \\ &= \argmin_{b \in \mathcal{B}} \sum_x P_t^{(u)}(x) l(b, x) \\ &= \argmin_{b \in \mathcal{B}} \left[ c_t \cdot \sum_x N_{t-1}^{(p)}(x) l(b, x) \right] \\ &= \mathbf{B}\{L_{t-1}^{(p)}(b)\}. \end{aligned} \quad (14)$$

Notice that the constant  $c_t$  does not affect the minimum.

2) *Bounding the expected loss:* In terms of the expected loss  $\mathbb{E} [\hat{L}_n] = \sum_{t=1}^n \mathbb{E} [l(\hat{b}_t, x_t)]$ , only the marginal distribution of  $\hat{b}_t$  matters, and therefore dependence between  $u_t(x)$  at different times does not affect the expected loss. Therefore in this section, we assume all  $u_t$  are equal,  $u_t(x) = u_1(x)$ .

We start by analyzing a clairvoyant predictor which includes also the state  $x_t$  into the prediction. For this purpose, let us define analogously to (12)-(13):

$$N_t^{(\text{cl})}(x) \triangleq N_t(x) + h_t u_t(x), \quad L_t^{(\text{cl})}(b) \triangleq \sum_x l(b, x) N_t^{(\text{cl})}(x) \quad (15)$$

where for  $t=0$  we define  $h_0 = 0$ , and note that  $N_0(x) = 0$  by definition, and therefore  $N_0^{(\text{cl})}(x) = 0$  and  $L_0^{(\text{cl})}(b) = 0$ . We consider the loss of the predictor  $\hat{b}_t = \mathbf{B}\{L_t^{(\text{cl})}(b)\}$ :

$$\begin{aligned} &\sum_{t=1}^n l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x_t) \\ &\stackrel{(a)}{=} \sum_{t=1}^n \sum_x [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x) (N_t(x) - N_{t-1}(x))] \\ &= \sum_{t=1}^n \sum_x [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x) (N_t^{(\text{cl})}(x) - N_{t-1}^{(\text{cl})}(x))] \\ &\quad - \sum_{t=1}^n \sum_x [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x) (h_t u_t(x) - h_{t-1} u_{t-1}(x))] \end{aligned} \quad (16)$$

where in (a) we used  $N_t(x) - N_{t-1}(x)$  as an indicator function

$\text{Ind}(x_t = x)$ . The first part can be bounded as:

$$\begin{aligned}
& \sum_{t=1}^n \sum_x [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x)(N_t^{(\text{cl})}(x) - N_{t-1}^{(\text{cl})}(x))] \\
&= \sum_{t=1}^n [L_t^{(\text{cl})}(\mathbf{B}\{L_t^{(\text{cl})}(b)\}) - L_{t-1}^{(\text{cl})}(\mathbf{B}\{L_t^{(\text{cl})}(b)\})] \\
&\stackrel{(a)}{\leq} \sum_{t=1}^n [L_t^{(\text{cl})}(\mathbf{B}\{L_t^{(\text{cl})}(b)\}) - L_{t-1}^{(\text{cl})}(\mathbf{B}\{L_{t-1}^{(\text{cl})}(b)\})] \\
&\stackrel{(b)}{=} L_n^{(\text{cl})}(\mathbf{B}\{L_n^{(\text{cl})}(b)\}) - L_0^{(\text{cl})}(\mathbf{B}\{L_0^{(\text{cl})}(b)\}) \\
&= L_n^{(\text{cl})}(\mathbf{B}\{L_n^{(\text{cl})}(b)\}) \\
&\leq L_n^{(\text{cl})}(\mathbf{B}\{L_n(b)\}) \\
&= L_n(\mathbf{B}\{L_n(b)\}) + h_n \sum_x l(\mathbf{B}\{L_n(b)\}, x) u_n(x) \\
&\leq L_n(\mathbf{B}\{L_n(b)\}) + R|\mathcal{X}|h_n \\
&= L_n^* + R|\mathcal{X}|h_n,
\end{aligned} \tag{17}$$

where we used (a) the fact that  $\mathbf{B}\{L_{t-1}^{(\text{cl})}(b)\}$  is optimized for  $L_{t-1}^{(\text{cl})}$  and (b) the sum of the telescopic series. For the second sum in (16), let us use the assumption  $u_t = u_1$ . Then:

$$\begin{aligned}
& \left| \sum_{t=1}^n \sum_x [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x)(h_t u_t(x) - h_{t-1} u_{t-1}(x))] \right| \\
&= \left| \sum_{t=1}^n \sum_x [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x)(h_t - h_{t-1})u_1(x)] \right| \\
&\leq \sum_{t=1}^n \sum_x R|h_t - h_{t-1}| \\
&\stackrel{(a)}{=} |\mathcal{X}|R \sum_{t=1}^n (h_t - h_{t-1}) \\
&\stackrel{(b)}{=} R|\mathcal{X}|h_n,
\end{aligned} \tag{18}$$

where we used (a) the assumption that the sequence  $h_t$  is non decreasing, and (b) the definition  $h_0 = 0$ . Combining (18) and (17) into (16) yields:

$$\sum_{t=1}^n l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x_t) \leq L_n^* + 2R|\mathcal{X}|h_n \tag{19}$$

The next step is to bound the performance difference between the clairvoyant predictor  $\mathbf{B}\{L_t^{(\text{cl})}(b)\}$  and the FPF forecaster  $b_t^{(\text{FPF})} = \mathbf{B}\{L_{t-1}^{(\text{p})}(b)\}$ . The key is that the new element added to  $L_t^{(\text{cl})}(b)$  is  $l(b, x_t)$ , and the dither element  $u(x_t)$  (i.e. belonging to the state that actually happened at time  $t$ ) contributes an offset in the same direction, which cancels this addition or most values of  $u(x_t)$ . For this purpose let us write the accumulated losses as:

$$\begin{aligned}
L_{t-1}^{(\text{p})}(b) &= \sum_x l(b, x)(N_{t-1}(x) + h_t u_t(x)) \\
&= L_c + h_t u(x_t) l(b, x_t) \\
L_t^{(\text{cl})}(b) &= \sum_x l(b, x)(N_t(x) + h_t u_t(x)) \\
&= L_{t-1}^{(\text{p})}(b) + l(b, x_t) \\
&= L_c + (h_t u(x_t) + 1) l(b, x_t)
\end{aligned} \tag{20}$$

where we defined

$$L_c = \sum_x l(b, x) N_{t-1}(x) + \sum_{x \neq x_t} l(b, x) h_t u_t(x). \tag{21}$$

Noticing that the common part  $L_c$  is independent of  $u(x_t)$ , we compute the conditional expectation given  $L_c$  for each of the predictors:

$$\begin{aligned}
& \mathbb{E} [l(\mathbf{B}\{L_{t-1}^{(\text{p})}(b)\}, x_t) | L_c] \\
&= \mathbb{E} [l(\mathbf{B}\{L_c + h_t u(x_t) l(b, x_t)\}, x_t) | L_c] \\
&= \int_{v=0}^1 l(\mathbf{B}\{L_c + h_t l(b, x_t) v\}, x_t) dv \\
&= \int_{v=0}^1 g(v) dv,
\end{aligned} \tag{22}$$

where we defined for brevity  $g(v) = l(\mathbf{B}\{L_c + h_t l(b, x_t) v\}, x_t)$ , and

$$\begin{aligned}
& \mathbb{E} [l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x_t) | L_c] \\
&= \mathbb{E} [l(\mathbf{B}\{L_c + (h_t u(x_t) + 1) l(b, x_t)\}, x_t) | L_c] \\
&= \int_{v=0}^1 l(\mathbf{B}\{L_c + (h_t v + 1) l(b, x_t)\}, x_t) dv \\
&= \int_{v=0}^1 l(\mathbf{B}\{L_c + h_t (v + h_t^{-1}) l(b, x_t)\}, x_t) dv \\
&= \int_{v=h_t^{-1}}^{1+h_t^{-1}} l(\mathbf{B}\{L_c + h_t l(b, x_t) v\}, x_t) dv \\
&= \int_{v=h_t^{-1}}^{1+h_t^{-1}} g(v) dv
\end{aligned} \tag{23}$$

The integrands in (22),(23) are equal. Let us temporarily assume that for  $t \geq 1$ ,  $h_t > 1$ , so that the integration regions partially overlap. For most of the integration region, because the integrands are the same (no matter what  $l(\cdot, x_t)$  evaluates to), and the integration regions overlap, they cancel out, and we remain with the contribution of the edges where there is no overlap:

$$\begin{aligned}
& \mathbb{E} [l(\mathbf{B}\{L_{t-1}^{(\text{p})}(b)\}, x_t) - l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x_t) | L_c] \\
&\stackrel{(22),(23)}{=} \int_{v=0}^1 g(v) dv - \int_{v=h_t^{-1}}^{1+h_t^{-1}} g(v) dv \\
&= \int_0^{h_t^{-1}} g(v) dv - \int_1^{1+h_t^{-1}} g(v) dv \\
&\leq \int_{[0, h_t^{-1}] \cup [1, 1+h_t^{-1}]} |g(v)| dv \\
&\leq 2 \cdot R \cdot h_t^{-1}.
\end{aligned} \tag{24}$$

Recall that we assumed  $h_t \geq 1$ . For  $h_t \leq 1$  the bound (24) is trivially true (because the RHS is at least  $2R$ ), and therefore it holds for all  $h_t$ .

Applying the iterated expectations law and accumulating (24) yields:

$$\mathbb{E} \left[ \sum_{t=1}^n l(\mathbf{B}\{L_{t-1}^{(\text{p})}(b)\}, x_t) - \sum_{t=1}^n l(\mathbf{B}\{L_t^{(\text{cl})}(b)\}, x_t) \right] \leq 2R \sum_{t=1}^n h_t^{-1} \tag{25}$$

which, together with (19) yields:

$$\begin{aligned}
& \mathbb{E} \left[ \sum_{t=1}^n l(\mathbf{B}\{L_{t-1}^{(p)}(b)\}, x_t) \right] - L_n^* \\
&= \mathbb{E} \left[ \sum_{t=1}^n l(\mathbf{B}\{L_{t-1}^{(p)}(b)\}, x_t) - \sum_{t=1}^n l(\mathbf{B}\{L_t^{(cl)}(b)\}, x_t) \right] \\
&\quad + \mathbb{E} \left[ \sum_{t=1}^n l(\mathbf{B}\{L_t^{(cl)}(b)\}, x_t) \right] - L_n^* \\
&\stackrel{(25),(19)}{\leq} 2R \left( \sum_{t=1}^n h_t^{-1} + |\mathcal{X}|h_n \right) \triangleq \Delta,
\end{aligned} \tag{26}$$

which proves the first claim of the theorem.

3) *Choices of the dither amplitude sequence:* There remains the question of selecting the sequence of dither amplitudes  $h_t$ . For a given horizon  $n$ , a simple calculation shows, that the best choice in terms of minimizing  $\Delta$  is a constant  $h_t$ , which equals  $\sqrt{\frac{n}{|\mathcal{X}|}}$ . As will be seen below, a good choice of a varying  $h_t$  that yields an infinite horizon solution (i.e. in which the setting of  $h_t$  does not depend on the horizon  $n$ ) is  $h_t = \sqrt{\frac{2t}{|\mathcal{X}|}}$ . However, because in real life we do not choose an “optimal”  $h_t$ , it is first desired to show that for a wide range of choices, the resulting predictor’s expected regret tends to zero. This analysis is rather straightforward and reoccurs in many developments of this kind [1, Ex 4.7][10]. So, let us choose  $h_t = h_1 \cdot t^\alpha$ , with  $\alpha \in (0, 1)$ . Then,

$$\begin{aligned}
\sum_{t=1}^n h_t^{-1} &= h_1^{-1} \sum_{t=1}^n t^{-\alpha} \leq h_1^{-1} \left( 1 + \int_{x=1}^n x^{-\alpha} dx \right) \\
&= h_1^{-1} \left( 1 + \frac{1}{1-\alpha} (n^{1-\alpha} - 1) \right) \leq h_1^{-1} \frac{1}{1-\alpha} n^{1-\alpha}
\end{aligned} \tag{27}$$

Substituting in (26) yields:

$$\begin{aligned}
\frac{\Delta}{n} &= \frac{2R}{n} \left( \sum_{t=1}^n h_t^{-1} + |\mathcal{X}| \cdot h_n \right) \\
&\leq 2R \left( h_1^{-1} \frac{1}{1-\alpha} n^{-\alpha} + |\mathcal{X}| \cdot h_1 n^{\alpha-1} \right)
\end{aligned} \tag{28}$$

For any  $\alpha \in (0, 1)$  and any  $h_1$ , this yields  $\frac{\Delta}{n} \xrightarrow{n \rightarrow \infty} 0$ , i.e. Hannan’s consistency. It is straightforward to see that the best choice is obtained by  $\alpha = \frac{1}{2}$  and  $h_1 = \sqrt{\frac{2}{|\mathcal{X}|}}$ , which yields:

$$\frac{\Delta}{n} = \frac{2R}{\sqrt{n}} (2h_1^{-1} + |\mathcal{X}| \cdot h_1) = 4R \sqrt{\frac{2|\mathcal{X}|}{n}} \tag{29}$$

4) *Proof of Corollary 2.1:* To prove the corollary is it sufficient to notice that all strategies for which the loss is computed in the proof of Theorem 2, are in the aforementioned set of optimizing strategies. For the  $L_2$  loss it is easy to see that the set of optimizing strategies is the convex hull of  $\mathcal{X}$  (the strategy for given  $\lambda(x)$  can be interpreted as a center of mass of  $\mathcal{X}$  with varying weights to the different points).

## B. Proof of Theorem 3 (Log loss)

The sequence is  $x_t, t = 1, \dots, n$ . The accumulated loss for probability  $q_t$  is  $\sum_{t=1}^n \log \frac{1}{q_t(x_t)}$ . The best fixed  $q_t$  in hindsight is  $q_t(x) = \hat{P}_x(x)$  and yields  $L_n^* = n \sum_x \hat{P}_x(x) \log \frac{1}{\hat{P}_x(x)}$ . For the universal estimator proposed:  $P_t^{(u)}(x) = c_t^{-1} \cdot (N_{t-1}(x) + h_t u_t(x))$ , and it is easy to see that given  $P_t^{(u)}(x)$ , the choice of  $q_t(\cdot)$ , the probability distribution for the next state is just  $q_t(\cdot) = \operatorname{argmin}_q \mathbb{E} \log \frac{1}{q(X)} = P_t^{(u)}(x)$

$$\begin{aligned}
\mathbb{E}[\hat{L}_n] &= \mathbb{E} \sum_{t=1}^n \log \frac{1}{P_t^{(u)}(x_t)} \\
&= s \sum_{t=1}^n [\mathbb{E} \log(c_t) - \mathbb{E} \log(N_{t-1}(x_t) + h_t u_t(x_t))]
\end{aligned} \tag{30}$$

In general, in order to achieve a small regret for the log loss, it is required that the overall contribution of  $P_t^{(u)}(x_t)$  for all occurrences of a certain state  $x_t = x$ , would approximate  $\hat{P}_x(x)$ . However the most important property, which is not satisfied by FL, is not to give a probability too close to 0 for a certain state  $x$  on its first appearance in the sequence. I.e. if  $N_{t-1}(x_t) = 0$ , it is required that  $\mathbb{E} \log(N_{t-1}(x_t) + h_t u_t(x_t)) = \mathbb{E} \log(h_t u_t(x_t))$  is finite. Indeed it is easy to verify that this holds.

Following is the detailed calculation and bounding for  $\mathbb{E}[\hat{L}_n]$ . The normalized  $c_t$  is bounded as:

$$c_t = \sum_x (N_{t-1}(x_t) + h_t u_t(x)) \leq t - 1 + |\mathcal{X}|h_t \tag{31}$$

In the below, denote for conciseness  $N_{t-1}(x_t) = v$

$$\begin{aligned}
& \mathbb{E} \log(N_{t-1}(x_t) + h_t u_t(x_t)) \\
&= \mathbb{E} \log(v + h_t u_t(x_t)) = \int_0^1 \log(v + h_t y) dy \\
&= h_t^{-1} \int_v^{v+h_t} \log(y) dy \\
&= h_t^{-1} [y \log y - y]_v^{v+h_t} \\
&= h_t^{-1} [(v + h_t) \log(v + h_t) - v \log v - h_t] \\
&\stackrel{(a)}{=} h_t^{-1} [h_t \log(v + h_t) + v \log(1 + \frac{h_t}{v}) - h_t] \\
&\stackrel{\log x \geq \frac{x}{1+x}}{\geq} h_t^{-1} \left[ h_t \log(v + h_t) + v \frac{\frac{h_t}{v}}{1 + \frac{h_t}{v}} - h_t \right] \\
&= \log(v + h_t) - \frac{h_t}{v + h_t} \\
&= \log(N_{t-1}(x_t) + h_t) - \frac{h_t}{N_{t-1}(x_t) + h_t}.
\end{aligned} \tag{32}$$

In (a) notice that  $v = 0$  is a special case. using  $0 \log 0 = 0$ , it is easy to verify that in this case the the expression before (a) for  $v = 0$  equals  $\log(h_t) - 1$ , and therefore the inequality holds.



Returning to (30):

$$\begin{aligned}
\mathbb{E}[\hat{L}_n] &= \sum_{t=1}^n [\mathbb{E} \log(c_t) - \mathbb{E} \log(N_{t-1}(x_t) + h_t u_t(x))] \\
&\leq \sum_{t=1}^n \log(t-1 + |\mathcal{X}|h_t) \\
&\quad + \sum_{t=1}^n \left[ -\log(N_{t-1}(x_t) + h_t) + \frac{h_t}{N_{t-1}(x_t) + h_t} \right] \\
&= \sum_{t=1}^n \log(t-1 + |\mathcal{X}|h_t) \\
&\quad + \sum_{x \in \mathcal{X}} \sum_{t: x_t=x} \left[ -\log(N_{t-1}(x) + h_t) + \frac{h_t}{N_{t-1}(x) + h_t} \right] \tag{33}
\end{aligned}$$

For the second sum, which was broken into the subsequences in which a specific state  $x$  appears, notice that  $N_{t-1}(x)$  increases by 1 between consecutive elements of the internal sum. The final value of  $N_{t-1}(x)$  in the last element equals the total number of appearances of  $x$ ,  $N_n(x)$ . At this point, it is beneficial to write  $L_n^*$  in a similar form:

$$\begin{aligned}
L_n^* &= n \sum_x \hat{P}_x(x) \log \frac{1}{\hat{P}_x(x)} \\
&= \sum_x N_n(x) \log \frac{n}{N_n(x)} \tag{34} \\
&= n \log n - \sum_x N_n(x) \log N_n(x)
\end{aligned}$$

A consequence of the bound on the size of a type class  $|\mathcal{T}_P| \leq \exp(nH(P))$  [14, Lemma II.2] is (considering the type defined by the sequence  $x, (\dots, \frac{N_n(x)}{n}, \dots)$  and taking the log of both sides):

$$\begin{aligned}
\log \left( \frac{n!}{\prod_{x \in \mathcal{X}} N_n(x)!} \right) &\leq n \sum_{x \in \mathcal{X}} \frac{N_n(x)}{n} \log \frac{n}{N_n(x)} \\
&= n \log n - \sum_{x \in \mathcal{X}} N_n(x) \log N_n(x) \tag{35}
\end{aligned}$$

Plugging into (34) yields:

$$\begin{aligned}
L_n^* &\geq \log \left( \frac{n!}{\prod_{x \in \mathcal{X}} N_n(x)!} \right) \\
&= \sum_{t=1}^n \log(t) - \sum_{x \in \mathcal{X}} \sum_{m=1}^{N_n(x)} \log(m) \tag{36}
\end{aligned}$$

Notice that this way of bounding  $L_n^*$  is slightly non standard: rather than writing  $L_n^*$  in a similar form to  $L_U$ , it would generally be simpler to write the bound on  $\mathbb{E}[\hat{L}_n]$  using factorials, and simplify it using Stirling's approximation, obtaining a form similar to (34), however this approach does not hold for varying  $h_t$ .

Let us now assume  $h_t$  is non-decreasing. Combining (33) with (36) yields:

$$\begin{aligned}
&\mathbb{E}[\hat{L}_n] - L_n^* \\
&\leq \sum_{t=1}^n (\log(t-1 + |\mathcal{X}|h_t) - \log(t)) \\
&\quad + \sum_{x \in \mathcal{X}} \sum_{t: x_t=x} \left[ \log(N_{t-1}(x) + 1) \right. \\
&\quad \left. - \log(N_{t-1}(x) + h_t) + \frac{h_t}{N_{t-1}(x) + h_t} \right] \\
&= \sum_{t=1}^n \left( \log \left( 1 + \frac{|\mathcal{X}|h_t - 1}{t} \right) \right) \\
&\quad + \sum_{x \in \mathcal{X}} \sum_{t: x_t=x} \left[ \log \left( 1 + \frac{1 - h_t}{N_{t-1}(x) + h_t} \right) \right. \\
&\quad \left. + \frac{h_t}{N_{t-1}(x) + h_t} \right] \tag{37} \\
&\leq \sum_{t=1}^n \frac{|\mathcal{X}|h_t - 1}{t} \\
&\quad + \sum_{x \in \mathcal{X}} \sum_{t: x_t=x} \left[ \frac{1 - h_t}{N_{t-1}(x) + h_t} + \frac{h_t}{N_{t-1}(x) + h_t} \right] \\
&= \sum_{t=1}^n \frac{|\mathcal{X}|h_t - 1}{t} + \sum_{x \in \mathcal{X}} \sum_{t: x_t=x} \frac{1}{N_{t-1}(x) + h_t} \\
&= \sum_{t=1}^n \frac{|\mathcal{X}|h_t - 1}{t} + \sum_{t=1}^n \frac{1}{N_{t-1}(x_t) + h_t}
\end{aligned}$$

Let us consider the sequence  $\mathbf{x}$  that maximizes the second sum. As clear intuitively, and will be proven below, this sequence selects all states of  $x$  in a round-robin fashion, which minimizes the growth rate of  $N_{t-1}(x_t)$  and for which  $N_{t-1}(x_t) = \lfloor \frac{t-1}{|\mathcal{X}|} \rfloor$ .

First, for a given type (i.e. for given  $\{N_n(x)\}_{x \in \mathcal{X}}$ ), consider the order that would yield the maximum. It is clear that the  $m$ -th occurrences of different states (assuming these states indeed occur at least  $m$  times) should occur at consecutive  $t$ -s. In other words, the sequence  $N_{t-1}(x_t)$  is non decreasing. Suppose that the opposite occurs, i.e.  $N_{t-2}(x_{t-1}) > N_{t-1}(x_t)$ , then obviously  $x_t \neq x_{t-1}$ . Let us flip the order of these states, i.e. let  $x'_t = x_{t-1}$  and  $x'_{t-1} = x_t$ , then as a result the counts will also flip,  $N'_{t-2}(x'_{t-1}) = N_{t-1}(x_t)$  and  $N'_{t-1}(x'_t) = N_{t-2}(x_{t-1})$ . This is easiest to see via an example: suppose the sequence is  $\mathbf{x} = (c, a, a, b, c, c)$ , then the counts  $N_{t-1}(x_t)$  are  $0, 0, 1, 0, 1, 2$ . After flipping the states  $t = 3, 4$  the sequence is  $\mathbf{x}' = (c, a, b, a, c, c)$  and the counts are  $0, 0, 0, 1, 1, 2$ . The elements pertaining to other times are not affected by this flip, while the sum of the two elements is now:

$$\begin{aligned}
&\frac{1}{N'_{t-2}(x'_{t-1}) + h_{t-1}} + \frac{1}{N'_{t-1}(x'_t) + h_t} \\
&= \frac{1}{N_{t-1}(x_t) + h_{t-1}} + \frac{1}{N_{t-2}(x_{t-1}) + h_t} \tag{38} \\
&> \frac{1}{N_{t-2}(x_{t-1}) + h_{t-1}} + \frac{1}{N_{t-1}(x_t) + h_t}
\end{aligned}$$

where the inequality holds because  $h_t \geq h_{t-1}$  and  $N_{t-2}(x_{t-1}) > N_{t-1}(x_t)$ . This can be seen by direct algebraic



manipulation, or by using Lemma 1 with the convex function  $f(x) = \frac{1}{x}$ :

**Lemma 1.** *Let  $f$  be a convex- $\cup$  function and  $a_0, a_1, b_0, b_1$  satisfy  $a_1 \geq a_0$  and  $b_1 \geq b_0$ , then*

$$f(a_0 + b_0) + f(a_1 + b_1) \geq f(a_0 + b_1) + f(a_1 + b_0) \quad (39)$$

*I.e. the maximum sum is obtained by joining the smaller and bigger elements together.*

*Proof:* Let us assume there is strict inequality at least in one of the pairs, otherwise the result holds trivially. Write the hybrid sums as convex combinations of the homogenous sums:

$$\begin{aligned} a_0 + b_1 &= \lambda(a_0 + b_0) + (1 - \lambda)(a_1 + b_1) \\ a_1 + b_0 &= (1 - \lambda)(a_0 + b_0) + \lambda(a_1 + b_1), \end{aligned} \quad (40)$$

with

$$\lambda = \frac{a_1 - a_0}{a_1 - a_0 + b_1 - b_0} \in [0, 1]. \quad (41)$$

From (40) and the convexity of  $f$ :

$$\begin{aligned} f(a_0 + b_1) &\leq \lambda f(a_0 + b_0) + (1 - \lambda)f(a_1 + b_1) \\ f(a_1 + b_0) &\leq (1 - \lambda)f(a_0 + b_0) + \lambda f(a_1 + b_1). \end{aligned} \quad (42)$$

Summing the two equations in (42) yields the desired result (39).

The conclusion is that the sequence  $N_{t-1}(x_t)$  is non-decreasing. Next, is it obvious that to increase the sum, all states in  $\mathbf{x}$  should be chosen approximately the same number of times. If at some point, a state  $x$  has been chosen for the  $m$ -th time at time  $t$ , while another state  $x'$  did not appear  $m - 2$  times at this point, then because of the monotonicity of  $N_{t-1}(x_t)$ ,  $x'$  can never appear again in an optimal sequence. Clearly, choosing  $x'$  instead of  $x$  at time  $t$  would decrease  $N_{t-1}(x_t)$  for time  $t$  as well as for all future occurrences of the state  $x$ , and therefore will increase the sum.

This concludes the proof that the last sum in (37) is maximized by  $N_{t-1}(x_t) = \lfloor \frac{t-1}{|\mathcal{X}|} \rfloor$ . Now, (37) may be rewritten as:

$$\mathbb{E}[\hat{L}_n] - L_n^* \leq \sum_{t=1}^n \frac{|\mathcal{X}|h_t - 1}{t} + \sum_{t=1}^n \frac{1}{\lfloor \frac{t-1}{|\mathcal{X}|} \rfloor + h_t} \triangleq \Delta. \quad (43)$$

The last bound is only a function of  $\{h_t\}$  and not of the sequence  $\mathbf{x}$ . If  $h_t$  grows sublinearly, then the dominant factor in the second sum will be  $\lfloor \frac{t-1}{|\mathcal{X}|} \rfloor$  and the sum would grow like  $O(\log(n))$ . On the other hand, any growth rate of  $h_t$  that satisfies  $\frac{1}{n} \sum_{t=1}^n \frac{h_t}{t} \xrightarrow{n \rightarrow \infty} 0$  would yield normalized expected regret tending to 0, and a sufficient condition is  $\frac{h_t}{t} \rightarrow 0$ , and particularly this holds for  $h_t = O(t^\alpha)$ ,  $\alpha \in [0, 1)$ .

If  $h_t$  is constant, then the first sum grows like  $O(\log(n))$

as well. A more detailed evaluation yields:

$$\begin{aligned} \Delta &= (|\mathcal{X}|h - 1) \sum_{t=1}^n \frac{1}{t} + \sum_{t=1}^n \frac{1}{\lfloor \frac{t-1}{|\mathcal{X}|} \rfloor + h} \\ &\leq (|\mathcal{X}|h - 1) \left( 1 + \int_1^n \frac{1}{t} dt \right) + \int_{t=0}^n \frac{1}{\lfloor \frac{t}{|\mathcal{X}|} \rfloor - 1 + h} dt \\ &= (|\mathcal{X}|h - 1) (1 + \log(n)) + |\mathcal{X}| \log \left( \frac{n - |\mathcal{X}|}{|\mathcal{X}|h} + 1 \right) \\ &\stackrel{h=|\mathcal{X}|^{-1}}{=} |\mathcal{X}| \log(n - |\mathcal{X}| + 1) \\ &\leq |\mathcal{X}| \log(n). \end{aligned} \quad (44)$$

This ends the proof of Theorem 3.

### C. Proof of Theorem 1

Both Theorem 2 and Theorem 3 show the normalized expected regret tends to 0 with  $n$ , and it remains to change from claims on expected regret to claims on the almost-sure regret. To prove that the regret tends to 0 almost surely, or more precisely,  $\limsup_{n \rightarrow \infty} (\hat{L}_n - L_n^*) \leq 0$ , using the already established fact that  $\limsup_{n \rightarrow \infty} (\mathbb{E}[\hat{L}_n] - L_n^*) \leq 0$ , it remains to show that  $\frac{1}{n} (\hat{L}_n - \mathbb{E}[\hat{L}_n]) \xrightarrow[n \rightarrow \infty]{} 0$  almost surely.

$$\frac{1}{n} (\hat{L}_n - \mathbb{E}[\hat{L}_n]) = \frac{1}{n} \sum_{t=1}^n \left( l(\hat{b}_t, x_t) - \mathbb{E}[l(\hat{b}_t, x_t)] \right) \quad (45)$$

To show that the mean above converges to 0 almost surely, we use Kolmogorov's criterion for the applicability of the Strong Law of Large numbers [15]. The elements of the sequence  $\gamma_t = l(\hat{b}_t, x_t) - \mathbb{E}[l(\hat{b}_t, x_t)]$  have zero mean, and are independent, because each  $\hat{b}_t$  depends only on the deterministic history of the sequence, and on  $u_t(x)$ , which are assumed independent. Notice that  $\gamma_t$  are not identically distributed. Kolmogorov's criterion requires that  $\sum_{t=1}^{\infty} \frac{\text{Var}(\gamma_t)}{t^2} < \infty$ . This holds trivially for bounded loss functions, for which the boundness of  $|\gamma_t|$  yields a constant bound on its variance. Proving that this condition holds for the case of the log loss function is a rather technical calculation which is deferred to Appendix-A.  $\square$

## APPENDIX

### A. Completion of the proof of Theorem 1 for the log loss

This appendix completes the proof of Theorem 1 from Section V-C, by showing that for the case of the log loss, the Kolmogorov criterion holds and therefore the normalized regret converges almost surely to the normalized expected regret. Our purpose is to upper bound the following variance:

$$\sigma_t^2 = \text{Var}(\gamma_t) = \text{Var}(l(\hat{b}_t, x_t)) = \text{Var} \left[ \log \left( P_t^{(u)}(x_t) \right) \right], \quad (46)$$

and show that Kolmogorov's criterion  $\sum_{t=1}^{\infty} \frac{\sigma_t^2}{t^2} < \infty$  holds.  $P_t^{(u)}(x_t)$  is defined in (7) as

$$P_t^{(u)}(x) = \frac{N_{t-1}(x) + h_t \cdot u_t(x)}{t - 1 + h_t \cdot \sum_{x' \in \mathcal{X}} u_t(x')}. \quad (47)$$

We have:

$$\begin{aligned}
\log \left( P_t^{(u)}(x_t) \right) &= \log(N_{t-1}(x_t) + h_t \cdot u_t(x_t)) \\
&\quad - \log(t-1 + h_t \cdot \sum_{x' \in \mathcal{X}} u_t(x')) \\
&= \underbrace{\log \left( \frac{N_{t-1}(x_t)}{h_t} + u_t(x_t) \right)}_A \\
&\quad - \underbrace{\log \left( \frac{t-1}{h_t} + \sum_{x' \in \mathcal{X}} u_t(x') \right)}_B.
\end{aligned} \tag{48}$$

Using

$$\text{Var}(A - B) \leq E[(A - B)^2] \leq 2E[A^2] + 2E[B^2], \tag{49}$$

where the second inequality stems from  $(a - b)^2 = 2a^2 + 2b^2 - (a + b)^2 \leq 2a^2 + 2b^2$ , it is enough to bound the expected squared value of each log in (48) separately. For a uniform r.v.  $U \sim \mathbb{U}[0, 1]$  and a constant  $a \geq 0$ , the following bound holds:

$$\begin{aligned}
\mathbb{E}[\log^2(a + U)] &= \int_0^1 \log^2(a + u) du \\
&= \int_a^{a+1} \log^2(u) du \\
&\leq \int_0^1 \log^2(u) du + \int_{\max(a, 1)}^{a+1} \log^2(u) du \\
&\leq \int_0^1 \log^2(u) du + \int_{\max(a, 1)}^{a+1} \log^2(u) du \\
&\leq [u \log^2(u) - 2u \log(u) + u]_0^1 + \log^2(a + 1) \\
&= 2 + \log^2(a + 1).
\end{aligned} \tag{50}$$

We used the fact that  $\log^2(u)$  is increasing for  $u \geq 1$ . Notice that the bound is trivial for  $a \geq 1$  because  $U \leq 1$ . Applying the bound to the squared elements in (48):

$$\begin{aligned}
\mathbb{E} \left[ \log^2 \left( \frac{N_{t-1}(x_t)}{h_t} + u_t(x_t) \right) \right] &\leq 2 + \log^2 \left( \frac{N_{t-1}(x_t)}{h_t} + 1 \right) \\
&\leq 2 + \log^2 \left( \frac{t-1}{h_t} + 1 \right).
\end{aligned} \tag{51}$$

For the second term, the expectation is first applied only to one arbitrary element  $u_t(x)$ , while conditioning on the other elements:

$$\begin{aligned}
&\mathbb{E} \left[ \log^2 \left( \frac{t-1}{h_t} + \sum_{x' \in \mathcal{X}} u_t(x') \right) \right] \\
&= \mathbb{E} \left\{ \mathbb{E} \left[ \log^2 \left( \frac{t-1}{h_t} + \sum_{x' \in \mathcal{X}} u_t(x') \right) \middle| u_t(x'), x' \neq x \right] \right\} \\
&\stackrel{(50)}{\leq} \mathbb{E} \left\{ 2 + \log^2 \left( \frac{t-1}{h_t} + \sum_{x' \neq x} u_t(x') + 1 \right) \right\} \\
&\leq 2 + \log^2 \left( \frac{t-1}{h_t} + |\mathcal{X}| \right),
\end{aligned} \tag{52}$$

where we used again the fact that  $\log^2(u)$  is increasing for  $u \geq 1$ . Combining (48) with (49) and the bounds above yields:

$$\begin{aligned}
\sigma_t^2 &= \text{Var} \left[ \log \left( P_t^{(u)}(x_t) \right) \right] \\
&\leq 4 + 2 \log^2 \left( \frac{t-1}{h_t} + 1 \right) + 4 + 2 \log^2 \left( \frac{t-1}{h_t} + |\mathcal{X}| \right) \\
&\leq 8 + 4 \log^2 \left( \frac{t-1}{h_t} + |\mathcal{X}| \right)
\end{aligned} \tag{53}$$

Under the assumptions of Theorem 1,  $h_t = h_1 \cdot t^\alpha$  with  $\alpha \in (0, 1)$  and so  $\sigma_t^2 = O(\log^2(t))$  and clearly Kolmogorov's criterion holds.  $\square$

## REFERENCES

- [1] N. Cesa-Bianchi and G. Lugosi, *Prediction, learning and games*. Cambridge University Press, 2006.
- [2] N. Merhav and M. Feder, "Universal schemes for sequential decision from individual data sequences," *IEEE Trans. Information Theory*, vol. 39, no. 4, pp. 1280–1292, Jul. 1993.
- [3] —, "Universal prediction," *IEEE Trans. Information Theory*, vol. 44, no. 6, pp. 2124–2147, Oct. 1998.
- [4] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley & sons, 1991.
- [5] M. Feder, "Gambling using a finite state machine," *IEEE Trans. Information Theory*, vol. 37, no. 5, pp. 1459–1465, sep 1991.
- [6] R. E. Krichevsky and V. K. Trofimov, "The performance of universal encoding," *IEEE Trans. Information Theory*, vol. 27, no. 2, pp. 199–207, Mar. 1981.
- [7] J. Hannan, "Approximation to bayes risk in repeated play," *Princeton University Press*, vol. Contributions to the Theory of Games, III, Ann. Math. Study Number 39, pp. 97–139, 1957.
- [8] M. Feder, N. Merhav, and M. Gutman, "Universal prediction of individual sequences," *IEEE Trans. Information Theory*, vol. 38, no. 4, Jul. 1992.
- [9] T. M. Cover, "Behavior of sequential predictors of binary sequences," in *Proc. 4th Prague Con. Inform. Theory, Statistical Decision Functions, Random Processes*, 1965, pp. 263–272.
- [10] A. T. Kalai and S. Vempala, "Efficient algorithms for online decision problems," *Journal of Computer and System Sciences*, vol. 71, no. 3, pp. 291–307, Oct. 2005.
- [11] R. Weatherford, *Philosophical foundations of probability theory*. London : Routledge & Kegan Paul, 1982.
- [12] Z. Reznikova and B. Ryabko, "Ants and bits," *IEEE Information Theory Society Newsletter*, vol. 62, no. 5, pp. 17–20, 2012.
- [13] Y. Lomnitz, "Universal communication over unknown channels," Ph.D. dissertation, Tel Aviv University, Aug. 2012, available online [http://www.eng.tau.ac.il/~yuval/publications/Yuval\\_PhD\\_report.pdf](http://www.eng.tau.ac.il/~yuval/publications/Yuval_PhD_report.pdf).
- [14] I. Csiszár, "The method of types," *IEEE Trans. Information Theory*, vol. 44, no. 6, pp. 2505–2523, Oct. 1998.
- [15] W. Feller, *An Introduction to Probability Theory and Its Applications*. New York: Wiley, 1971.