# Minimax Online Prediction of Varying Bernoulli Process under Variational Approximation

**Kenta Konagayoshi**                    KONAGAYOSHI.KENTA@INF.KYUSHU-U.AC.JP
*Department of Informatics, Kyushu University, Fukuoka, 819-0395 Japan*

**Kazuho Watanabe**                    WKAZUHO@CS.TUT.AC.JP
*Department of Computer Science and Engineering, Toyohashi University of Technology, Toyohashi, 441-8580 Japan*

## Abstract

We consider the online prediction of a varying Bernoulli process (sequence of varying Bernoulli probabilities) from a single binary sequence. A real-valued online prediction method has been proposed as a prior work that incorporates the smoothness of the prediction sequence into the concept of the regret. Also, a Bayesian prediction method for the varying Bernoulli processes has been developed based on the variational inference. However, the former is not applicable to loss functions other than the squared error function, and the latter has no guarantee on the regret as an online prediction method. We propose a new online prediction method of a varying Bernoulli process from a single binary sequence with a guarantee to minimize the maximum regret under variational approximation. Through numerical experiments, we compare the Bayesian prediction method with the proposed method by using the regret with/without approximation and the KL divergence from the true underlying process. We discuss the prediction accuracy and influences of the approximation of the proposed method.

**Keywords:** online prediction, varying Bernoulli process, minimax strategy, regret, variational approximation

## 1. Introduction

Online prediction methods sequentially predict future data from time series data, and solve the concern of insufficient storage since it is unnecessary to keep historical data to obtain the prediction sequentially. Examples of applications include stock price prediction, weather forecast, object tracking, and so on. The regret is widely adopted as a measure of accuracy in online prediction. Since the regret expresses the difference between the cumulative losses of offline and online predictions, the smaller the regret, the closer the online prediction is to the optimal prediction. There have been developed online prediction methods that minimizes the maximum regret with some extensions (Cesa-Bianchi and Lugosi, 2006). Cutting edge methods include collaborative filtering bandits, sequential choice bandits, and so on (Li et al., 2015; Cao et al., 2019). In particular, some of online prediction methods take into account the nonstationary setting where the best offline predictor can change over time (Moroshko and Crammer, 2014). Herbster and Warmuth (2001) introduced the smoothness of the shift in the predictor under the nonstationary setting. Such a smoothness

has been directly modeled in the Bayesian framework by the Kalman filter (Kalman, 1960) or more generally by state space models. More recently, Koolen et al. (2015) proposed a prediction method that minimizes the maximum regret where the offline predictor is assumed to have the smoothness by directly incorporating a quadratic regularization term in a similar fashion to the Kalman filter. Although the authors succeeded in developing an efficient implementation of this method, the key ingredients for the efficient computation are the squared loss function and the quadratic regularization term (Koolen et al., 2014). It has been considered difficult to use other loss functions because the derivation of this method strongly depends on the squared loss function. This is analogous to the conjugacy of the Gaussian likelihood and prior in Bayesian methods where other likelihood nonconjugate to the Gaussian prior demands analytically intractable computation.

In this study, we consider the prediction of the varying sequence of probabilities that an event occurs, which we refer to as the varying Bernoulli process. More specifically, we predict a sequence of Bernoulli distributions with varying parameters. Such a model is included in the class of arbitrarily varying sources in information theory (Han and Kobayashi, 2007). The prediction of such a source is desirable for predictive coding schemes such as the arithmetic coding under nonstationary environments (Han and Kobayashi, 2007; Rissanen and Langdon, 1981). Similar models have also been applied in the field of neuroscience and communication engineering. For example, neuronal firing rate and congestion degree of channel are predicted from binary sequences representing the presence or absence of an event such as neuronal firing and arrival of communication packets. Bayesian prediction methods which predict the varying Bernoulli processes from a binary sequence were proposed (Koyama and Shinomoto, 2005; Cunningham et al., 2009; Watanabe and Okada, 2011; Takiyama and Okada, 2010). However, these prediction methods have no guarantee on the regret. Therefore, we propose a new online prediction method of varying Bernoulli processes from a single binary sequence. This is achieved by extending the framework of the previous study by Koolen et al. (2015) to the logistic loss function, and applying the variational approximation to it, which has been used for Bayesian inference in the logistic regression (Jaakkola and Jordan, 2000). This method ensures that the maximum regret over all binary sequences is minimized under the variational approximation. That is, by achieving the minimax regret, the worst case prediction can have the theoretical guarantee that it will not be worse than this. In the prediction of varying Bernoulli processes, even the optimal offline prediction is analytically intractable, and hence there is little hope for direct derivation of an efficient online prediction algorithm. One of the main contributions of this paper is to demonstrate that the variational approximation enables us to apply the framework of Koolen et al. (2015), which leads to an efficient online prediction algorithm for this problem too. Moreover, we theoretically evaluate the upper and lower bounds of the minimax regret, and show that it grows as $T/\sqrt{\lambda_T}$, in the same order as the case of the squared loss (Koolen et al., 2015), where $T$ is the time horizon and $\lambda_T$ is the regularization parameter for the smoothness, which can also grow with $T$. We also numerically examine the prediction accuracy by the regret with or without approximation and the KL divergence from the true underlying process, and discuss the prediction accuracy and influences of the approximation of the proposed method. This paper contains only Appendix A while the supplementary material also contains Appendices B–E.

## 2. Settings

### 2.1. Definition

We handle the time series data $\boldsymbol{x} = \{x_t\}_{t=1}^{T}$ consisting of binary data $x_t \in \{0, 1\}$ that takes 0 if an event does not occur or 1 if an event occurs on the $t$th trial. We assume that $x_t$ follows the Bernoulli distribution with parameter $\theta_t \in [0, 1]$,

$$p(x_t|\theta_t) = \theta_t^{x_t}(1 - \theta_t)^{1-x_t}. \tag{1}$$

The logit transformation

$$a(\theta) = \ln \frac{\theta}{1 - \theta}$$

yields the following model from Eq. (1),

$$p(x_t|a_t) = \exp\{x_t a_t - \ln(1 + e^{a_t})\}, \tag{2}$$

where $a_t = a(\theta_t)$. By this transformation, parameter $\theta_t \in [0, 1]$ is converted into $a_t \in (-\infty, +\infty)$. We define the logistic error function derived from Eq. (2) by the log-loss, $-\ln p(x_t|a_t)$,

$$E(x_t, a_t) = -x_t a_t + \ln(1 + e^{a_t}) \tag{3}$$

and the regularizer inducing the smoothness of time series $\boldsymbol{a} = \{a_t\}_{t=1}^{T}$,

$$\sum_{t=1}^{T+1} (a_t - a_{t-1})^2,$$

which corresponds to the Gaussian prior on $\boldsymbol{a}$,

$$p(\boldsymbol{a}) \propto \exp\left\{-\lambda \sum_{t=1}^{T+1} (a_t - a_{t-1})^2\right\},$$

where the prior distribution induces the smoothness of the prediction sequence and let $a_0 = a_{T+1} = 0$. Here, $\lambda$ works as the regularization coefficient which determines the smoothness of the sequence to be predicted.

### 2.2. Regret

In this study, we use the regret, which is widely known in the field of online prediction, as a performance measure of a prediction algorithm. The regret evaluates the relative loss of the prediction sequentially observing time series data (online prediction) compared to the loss of the prediction after observing all time series data (offline prediction). Let the predicted sequence by online prediction be $a_1, \ldots, a_T$, predicted sequence by offline prediction be $\hat{a}_1, \ldots, \hat{a}_T$ and the loss function on an input $x$ and a prediction $a$ be $l(x, a)$. The regret incorporating the smoothness of the predicted sequence $R$ is defined by the following equation,

$$R = \sum_{t=1}^{T} l(x_t, a_t) - \min_{\hat{a}_1, \ldots, \hat{a}_T} \left\{ \sum_{t=1}^{T} l(x_t, \hat{a}_t) + \lambda \sum_{t=1}^{T+1} (\hat{a}_t - \hat{a}_{t-1})^2 \right\}. \tag{4}$$

The previous study by Koolen et al. (2015) derived an online prediction algorithm by setting $l(x, a)$ to the squared error function $(a - x)^2$, and by solving the minimax problem,

$$\min_{a_1} \max_{x_1} \cdots \min_{a_T} \max_{x_T} R,$$

where both $x_t$ and $a_t$ are real values.[1] In this study, we approximately solve the minimax problem using binary input data $x_t \in \{0, 1\}$, value to be predicted $a_t \in (-\infty, +\infty)$ and the logistic loss in Eq. (3) as the loss function $l$.

Let $\hat{\boldsymbol{a}} = (\hat{a}_1, \ldots, \hat{a}_T)^{\mathrm{T}}$. The regret (4) is generally expressed as

$$R = \sum_{t=1}^{T} l(x_t, a_t) - \min_{\hat{a}_1, \ldots \hat{a}_T} \left\{ \sum_{t=1}^{T} l(x_t, \hat{a}_t) + \lambda \hat{\boldsymbol{a}}^{\mathrm{T}} \boldsymbol{K} \hat{\boldsymbol{a}} \right\},$$

where $\boldsymbol{K}$ is an arbitrary positive definite matrix. The regret in Eq. (4) corresponds to the case where

$$\boldsymbol{K} = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & & \ddots & & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}. \tag{5}$$

Although we use this case as a running example of this paper, the discussion before Section 3.5 also applies to the general positive definite $\boldsymbol{K}$, as well as some results on the case of the squared loss (Koolen et al., 2015).

## 3. Method

### 3.1. Derivation of the Loss Function $l$

We consider solving the minimax problem using Eq. (3) as the loss function $l$. However, since the second term in Eq. (3) involves the nonquadratic term with respect to $a$, it is intractable to solve the resulting minimax problem analytically. Therefore, as in the previous works on Bayesian inference in logistic regression models, we approximate the loss function $E$ to a quadratic function using the variational approximation (Jaakkola and Jordan, 2000; Watanabe and Okada, 2011). First, we define

$$f(a^2) = \ln \left( e^{\frac{\sqrt{a^2}}{2}} + e^{-\frac{\sqrt{a^2}}{2}} \right). \tag{6}$$

By differentiating with $a^2$, the following equation is obtained,

$$f'(a^2) = \frac{1}{e^{\frac{\sqrt{a^2}}{2}} + e^{-\frac{\sqrt{a^2}}{2}}} \cdot \frac{1}{4\sqrt{a^2}} \cdot \left( e^{\frac{\sqrt{a^2}}{2}} - e^{-\frac{\sqrt{a^2}}{2}} \right) = \frac{1}{4\sqrt{a^2}} \cdot \tanh \frac{\sqrt{a^2}}{2}.$$

Also $f(a^2)$ is a concave function with respect to $a^2$. Therefore, we approximate Eq. (6) by the first order approximation on $a^2$. Here, we define $f'(a^2) = \phi(a^2)$. In addition, we

---

1. More precisely, the formulation in (Koolen et al., 2015) solves this minimax problem under some additional constraints on inputs such as the boundedness of $x_t$.

introduce variational parameter $\xi$ and we obtain the following inequality

$$f(a^2) \leq f(\xi^2) + \phi(\xi^2)(a^2 - \xi^2).$$

Since $\ln(1 + e^a) = f(a^2) + \frac{a}{2}$, we have

$$E(x, a) \leq -xa + \frac{a}{2} + f(\xi^2) + \phi(\xi^2)(a^2 - \xi^2). \tag{7}$$

We define the loss function $l$ in Eq. (4), by the right hand side of Eq. (7), that is,

$$l(x_t, a_t) = -x_t a_t + \frac{a_t}{2} + f(\xi_t^2) + \phi(\xi_t^2)(a_t{}^2 - \xi_t{}^2). \tag{8}$$

Thus, the regret in Eq. (4) is obtained as

$$R = \sum_{t=1}^{T} \left\{ -x_t a_t + \frac{a_t}{2} + f(\xi_t^2) + \phi(\xi_t^2)(a_t^2 - \xi_t^2) \right\} \tag{9}$$

$$- \min_{\hat{a}_1, \ldots, \hat{a}_T} \left\{ \sum_{t=1}^{T} \left\{ -x_t \hat{a}_t + \frac{\hat{a}_t}{2} + f(\xi_t^2) + \phi(\xi_t^2)(\hat{a}_t^2 - \xi_t^2) \right\} + \lambda \sum_{t=1}^{T+1} (\hat{a}_t - \hat{a}_{t-1})^2 \right\}.$$

Note that the minimax solution of this regret, which will be described later in this section, is not directly derived from the minimax solution under the squared loss obtained in (Koolen et al., 2015).

## 3.2. Offline Prediction and its Loss

The second term in Eq. (9) shows the loss caused by offline prediction,

$$L = \sum_{t=1}^{T} \left\{ -x_t \hat{a}_t + \frac{\hat{a}_t}{2} + f(\xi_t^2) + \phi(\xi_t^2)(\hat{a}_t^2 - \xi_t^2) \right\} + \lambda \sum_{t=1}^{T+1} (\hat{a}_t - \hat{a}_{t-1})^2. \tag{10}$$

After observing all time series data $\boldsymbol{x} = (x_1, \ldots, x_T)^{\mathrm{T}}$, the sequence $\hat{\boldsymbol{a}} = (\hat{a}_1, \ldots, \hat{a}_T)^{\mathrm{T}}$ that minimizes $L$ is the optimal offline prediction. Since $L$ is quadratic function of $\hat{\boldsymbol{a}}$, we can explicitly express $\hat{\boldsymbol{a}}$ that minimizes $L$ with $\boldsymbol{x}$ and $\boldsymbol{\xi} = (\xi_1, \ldots, \xi_T)^{\mathrm{T}}$. To solve minimization problem, express $L$ with vectors and matrices. Defining $\boldsymbol{f} = (f(\xi_1^2), \ldots, f(\xi_T^2))^{\mathrm{T}}$,

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi(\xi_1^2) & & \\ & \ddots & \\ & & \phi(\xi_T^2) \end{pmatrix},$$

we can express $L$ as follows:

$$L = -\boldsymbol{x}^{\mathrm{T}} \hat{\boldsymbol{a}} + \frac{1}{2} \mathbf{1}^{\mathrm{T}} \hat{\boldsymbol{a}} + \mathbf{1}^{\mathrm{T}} \boldsymbol{f} + \hat{\boldsymbol{a}}^{\mathrm{T}} \boldsymbol{\Phi} \hat{\boldsymbol{a}} - \boldsymbol{\xi}^{\mathrm{T}} \boldsymbol{\Phi} \boldsymbol{\xi} + \lambda \hat{\boldsymbol{a}}^{\mathrm{T}} \boldsymbol{K} \hat{\boldsymbol{a}}, \tag{11}$$

where $\mathbf{1}$ is a vertical vector whose size is $T$ and whose elements are all 1. Using Eq. (11), we find the prediction vector minimizing $L$, $\hat{\boldsymbol{a}}$, and the optimal loss $L^*$. The standard technique for minimizing the quadratic loss yields

$$\hat{\boldsymbol{a}} = \frac{1}{2} (\boldsymbol{\Phi} + \lambda \boldsymbol{K})^{-1} \left( \boldsymbol{x} - \frac{1}{2} \mathbf{1} \right),$$

$$L^* = -\frac{1}{4} \left( \boldsymbol{x} - \frac{1}{2} \mathbf{1} \right)^{\mathrm{T}} (\boldsymbol{\Phi} + \lambda \boldsymbol{K})^{-1} \left( \boldsymbol{x} - \frac{1}{2} \mathbf{1} \right) + \mathbf{1}^{\mathrm{T}} \boldsymbol{f} - \boldsymbol{\xi}^{\mathrm{T}} \boldsymbol{\Phi} \boldsymbol{\xi}. \tag{12}$$

### 3.3. Minimax Loss at One Time Point

We derive the minimax loss at one time point, which is used in Section 3.4. For an input $x \in \{0, 1\}$ and a prediction $a \in (-\infty, +\infty)$ as a prediction, let

$$V(a, x) = -xa + \frac{a}{2} + \phi(\xi^2)a^2 + \alpha \left( x - \frac{1}{2} \right) \tag{13}$$

be the loss at one time point and consider

$$V^* = \min_a \max_x V(a, x)$$

as the minimax loss at one time point. Here, $\alpha$ is an arbitrary real value. It follows from Eq. (13) that

$$V(a, 0) = \frac{a}{2} + \phi(\xi^2)a^2 - \frac{\alpha}{2}, \tag{14}$$

$$V(a, 1) = -\frac{a}{2} + \phi(\xi^2)a^2 + \frac{\alpha}{2}. \tag{15}$$

Moreover, completing the squares in Eq. (14) and Eq. (15) with respect to $a$, we have

$$V(a, 0) = \phi(\xi^2) \left( a + \frac{1}{4\phi(\xi^2)} \right)^2 - \frac{1}{16\phi(\xi^2)} - \frac{\alpha}{2}, \tag{16}$$

$$V(a, 1) = \phi(\xi^2) \left( a - \frac{1}{4\phi(\xi^2)} \right)^2 - \frac{1}{16\phi(\xi^2)} + \frac{\alpha}{2}. \tag{17}$$

Because it holds that $\phi(\xi^2) \geq 0$, the expressions in Eq. (16) and Eq. (17) represent two parabolas opening upward as illustrated in Fig. 1. Equalizing Eq. (16) and Eq. (17), we
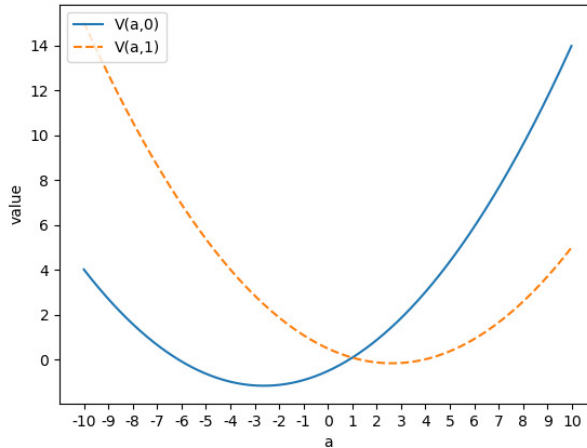


Figure 1: Two parabolas, Eq. (16) and Eq. (17).

find the optimal $a$ minimizing $\max_x V(a, x)$,

$$V(a, 0) = V(a, 1)$$
$$\Leftrightarrow \phi(\xi^2)\left(a + \frac{1}{4\phi(\xi^2)}\right)^2 - \frac{\alpha}{2} = \phi(\xi^2)\left(a - \frac{1}{4\phi(\xi^2)}\right)^2 + \frac{\alpha}{2}$$
$$\Leftrightarrow a = \alpha.$$

Therefore, the minimax optimal solution is $a = \alpha$, and the optimal minimax value $V^*$ is given by

$$V(\alpha, 0) = V(\alpha, 1) = \phi(\xi^2)\alpha^2. \tag{18}$$

### 3.4. Derivation of Online Prediction Method

The minimax regret is

$$R^* = \min_{a_1} \max_{x_1} \cdots \min_{a_T} \max_{x_T} R, \tag{19}$$

where $R$ is defined in Eq. (9). Using Eq. (19), solve the minimax problem. Let $\boldsymbol{x}_t = (x_1, \ldots, x_t)^\mathrm{T}$ and define $V(\boldsymbol{x}_T)$ and $V(\boldsymbol{x}_{t-1})$ by

$$V(\boldsymbol{x}_T) = -L^* \tag{20}$$

$$V(\boldsymbol{x}_{t-1}) = \min_{a_t} \max_{x_t}\left\{-x_t a_t + \frac{a_t}{2} + f(\xi_t^2) + \phi(\xi_t^2)(a_t^2 - \xi_t^2) + V(\boldsymbol{x}_t)\right\} \tag{21}$$

Then $R^*$ is

$$R^* = V(\boldsymbol{x}_0).$$

Using the decomposition

$$\boldsymbol{R}_t = \begin{pmatrix} \boldsymbol{A}_t & \boldsymbol{b}_t \\ \boldsymbol{b}_t^\mathrm{T} & c_t \end{pmatrix}, \tag{22}$$

we define a $(t-1) \times (t-1)$ matrix

$$\boldsymbol{R}_{t-1} = \boldsymbol{A}_t + \phi(\xi_t^2)\boldsymbol{b}_t \boldsymbol{b}_t^\mathrm{T}, \tag{23}$$

recursively from

$$\boldsymbol{R}_T = (\boldsymbol{\Phi}_T + \lambda \boldsymbol{K}_T)^{-1}. \tag{24}$$

The following theorem finally follows from Eq. (12) and Eq. (18).

**Theorem 1** *The minimax value $V(\boldsymbol{x}_t)$ and the optimal prediction strategy $a_t$ for the problem Eq. (19) are given by*

$$V(\boldsymbol{x}_t) = \frac{1}{4}\left(\boldsymbol{x}_t - \frac{1}{2}\mathbf{1}_t\right)^\mathrm{T} \boldsymbol{R}_t\left(\boldsymbol{x}_t - \frac{1}{2}\mathbf{1}_t\right) - \mathbf{1}_t^\mathrm{T} \boldsymbol{f}_t + \boldsymbol{\xi}_t^\mathrm{T} \boldsymbol{\Phi}_t \boldsymbol{\xi}_t + \frac{1}{16} \sum_{s=t+1}^{T} c_s, \tag{25}$$

$$a_t = \frac{1}{2}\left(\boldsymbol{x}_{t-1} - \frac{1}{2}\mathbf{1}_{t-1}\right)^\mathrm{T} \boldsymbol{b}_t.$$

147

The resulting minimax regret is given by

$$R^* = V(\boldsymbol{x}_0) = \frac{1}{16} \sum_{t=1}^{T} c_t.$$

**Proof** We use induction.

(i) In the case of $t = T$, $V(\boldsymbol{x}_T)$ is rewritten from Eq. (12) as follows:

$$V(\boldsymbol{x}_T) = \frac{1}{4} \left( \boldsymbol{x}_T - \frac{1}{2} \mathbf{1}_T \right)^{\mathrm{T}} (\boldsymbol{\Phi}_T + \lambda \boldsymbol{K}_T)^{-1} \left( \boldsymbol{x}_T - \frac{1}{2} \mathbf{1}_T \right) - \mathbf{1}_T^{\mathrm{T}} \boldsymbol{f}_T + \boldsymbol{\xi}_T^{\mathrm{T}} \boldsymbol{\Phi}_T \boldsymbol{\xi}_T$$
$$= - L^*,$$

which satisfies Eq. (20).

(ii) Assuming that $V(\boldsymbol{x}_t)$ is expressed by Eq. (25), we prove that $V(\boldsymbol{x}_{t-1})$ also satisfies Eq. (25). From Eqs. (21) and (22), we have

$$V(\boldsymbol{x}_{t-1}) = \min_{a_t} \max_{x_t} \left\{ -x_t a_t + \frac{a_t}{2} + f(\xi_t^2) + \phi(\xi_t^2)(a_t^2 - \xi_t^2) \right\}$$
$$+ \frac{1}{4} \left( \boldsymbol{x}_t - \frac{1}{2} \mathbf{1}_t \right)^{\mathrm{T}} \boldsymbol{R}_t \left( \boldsymbol{x}_t - \frac{1}{2} \mathbf{1}_t \right) - \mathbf{1}_t^{\mathrm{T}} \boldsymbol{f}_t + \boldsymbol{\xi}_t^{\mathrm{T}} \boldsymbol{\Phi}_t \boldsymbol{\xi}_t + \frac{1}{16} \sum_{s=t+1}^{T} c_s$$
$$= \frac{1}{4} \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right)^{\mathrm{T}} \boldsymbol{A}_t \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right)$$
$$- \mathbf{1}_{t-1}^{\mathrm{T}} \boldsymbol{f}_{t-1} + \boldsymbol{\xi}_{t-1}^{\mathrm{T}} \boldsymbol{\Phi}_{t-1} \boldsymbol{\xi}_{t-1} + \frac{1}{16} \sum_{s=t+1}^{T} c_s + C, \tag{26}$$

where $C$ is given by

$$C = \min_{a_t} \max_{x_t} -x_t a_t + \frac{a_t}{2} + \phi(\xi_t^2) a_t^2 + \frac{c_t}{4} \left( x_t - \frac{1}{2} \right)^2 + \frac{1}{2} \left( x_t - \frac{1}{2} \right) \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right) \boldsymbol{b}_t.$$

Since $\frac{c_t}{4}(x_t - \frac{1}{2})^2$ is $\frac{c_t}{16}$ regardless of $x_t = 0$ or 1, the above function of $a_t$ and $x_t$ is expressed as $V(a_t, x_t)$ in Eq. (13) with $\xi = \xi_t$ and $\alpha = \frac{1}{2} \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right) \boldsymbol{b}_t$. Hence, the argument in Section 3.3 yields that the optimal prediction $a_t$ and $C$ are given by

$$a_t = \frac{1}{2} \left( \boldsymbol{x}_t - \frac{1}{2} \mathbf{1}_t \right)^{\mathrm{T}} \boldsymbol{b}_t, \tag{27}$$
$$C = \frac{1}{4} \phi(\xi_t^2) \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right)^{\mathrm{T}} \boldsymbol{b}_t \boldsymbol{b}_t^{\mathrm{T}} \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right) + \frac{c_t}{16}. \tag{28}$$

Substituting Eq. (27) and Eq. (28) into Eq. (26), and using Eq. (23), we obtain

$$V(\boldsymbol{x}_{t-1}) = \frac{1}{4} \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right)^{\mathrm{T}} \boldsymbol{R}_{t-1} \left( \boldsymbol{x}_{t-1} - \frac{1}{2} \mathbf{1}_{t-1} \right) - \mathbf{1}_{t-1}^{\mathrm{T}} \boldsymbol{f}_{t-1} + \boldsymbol{\xi}_{t-1}^{\mathrm{T}} \boldsymbol{\Phi}_{t-1} \boldsymbol{\xi}_{t-1} + \frac{1}{16} \sum_{s=t}^{T} c_s.$$

Therefore, $V(\boldsymbol{x}_{t-1})$ satisfies Eq. (25) and the theorem follows from (i) and (ii) by induction. ∎

### 3.5. Simplification of Online Prediction Formula in Special Case

If we use $\boldsymbol{K}$ in Eq. (5) and fix $\xi_t$ to a constant $\xi$ independent of $t$, we can simplify the online prediction formula. Therefore, focusing on this case, we derive a prediction formula and evaluate the upper and lower bound of the regret. First, we show the formula for obtaining each element of Eq. (24).

**Lemma 1** *(Hu and O'Connell, 1996)*
*For* $\sinh x = \frac{e^x - e^{-x}}{2}$ *and* $\cosh x = \frac{e^x + e^{-x}}{2}$, *let* $\nu = \cosh^{-1}\left(1 + \frac{\phi(\xi^2)}{2\lambda}\right)$. *Then the $ij$th element of* $(\boldsymbol{\Phi}_T + \lambda\boldsymbol{K}_T)^{-1}$ *is*

$$(\boldsymbol{\Phi}_T + \lambda\boldsymbol{K}_T)^{-1}_{i,j} = \frac{\cosh(\nu(T+1-|i-j|)) - \cosh(\nu(T+1-i-j))}{2\lambda\sinh(\nu)\sinh((T+1)\nu)}.$$

Here, we define $\boldsymbol{z}_t$, $h_t$ and $h$ as follows:

$$\begin{aligned}
\boldsymbol{z}_t =& (\boldsymbol{\Phi}_t + \lambda\boldsymbol{K}_t)^{-1}\boldsymbol{e}_t, \\
h_t =& \boldsymbol{e}_t^{\mathrm{T}}(\boldsymbol{\Phi}_t + \lambda\boldsymbol{K}_t)^{-1}\boldsymbol{e}_t = \boldsymbol{e}_t^{\mathrm{T}}\boldsymbol{z}_t, \\
h =& \frac{2}{1 + \frac{2\lambda}{\phi(\xi^2)} + \sqrt{1 + \frac{4\lambda}{\phi(\xi^2)}}},
\end{aligned} \tag{29}$$

where $\boldsymbol{e}_t$ is a vertical vector whose size is $t$ and in which the $t$th element is 1 and the other elements are 0.

**Lemma 2** *It holds that*

$$h_t = \frac{1}{\phi(\xi^2)} \cdot \frac{1 - (\frac{\lambda}{\phi(\xi^2)}h)^{2t}}{1 - (\frac{\lambda}{\phi(\xi^2)}h)^{2t+2}} \cdot h.$$

*In addition,* $\lim\limits_{t \to \infty} h_t = \dfrac{h}{\phi(\xi^2)}$.

Next, using the formula for finding an inverse matrix of the block matrix,

$$\begin{pmatrix} \boldsymbol{A} & \boldsymbol{b} \\ \boldsymbol{b}^{\mathrm{T}} & c \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{A}^{-1} + \boldsymbol{A}^{-1}\boldsymbol{b}(c - \boldsymbol{b}^{\mathrm{T}}\boldsymbol{A}^{-1}\boldsymbol{b})^{-1}\boldsymbol{b}^{\mathrm{T}}\boldsymbol{A}^{-1} & -\boldsymbol{A}^{-1}\boldsymbol{b}(c - \boldsymbol{b}^{\mathrm{T}}\boldsymbol{A}^{-1}\boldsymbol{b})^{-1} \\ -(c - \boldsymbol{b}^{\mathrm{T}}\boldsymbol{A}^{-1}\boldsymbol{b})^{-1}\boldsymbol{b}^{\mathrm{T}}\boldsymbol{A}^{-1} & (c - \boldsymbol{b}^{\mathrm{T}}\boldsymbol{A}^{-1}\boldsymbol{b})^{-1} \end{pmatrix},$$

we obtain recursive formulas for $h_t$ and $\boldsymbol{Z}_t$.

**Lemma 3** *It holds that*

$$\begin{aligned}
h_t =& \frac{1}{\phi(\xi^2) + 2\lambda - \lambda^2 h_{t-1}}, \\
\boldsymbol{z}_t =& h_t \begin{pmatrix} \lambda\boldsymbol{z}_{t-1} \\ 1 \end{pmatrix}.
\end{aligned} \tag{30}$$

Using Lemmas 1, 2 and 3, we obtain the prediction formula $a_t$ at time $t$.

**Theorem 2** *It holds that $\boldsymbol{R}_t^{-1} = \boldsymbol{\Phi}_t + \lambda \boldsymbol{K}_t + \gamma_t \boldsymbol{e}_t \boldsymbol{e}_t^{\mathrm{T}}$, and that*

$$a_t = \frac{1}{2}\lambda c_t \left( \boldsymbol{x}_{t-1} - \frac{1}{2}\boldsymbol{1}_{t-1} \right)^{\mathrm{T}} \boldsymbol{z}_{t-1}, \tag{31}$$

*where $\gamma_t = \frac{1}{c_t} - \frac{1}{h_t}$ and $c_t$ satisfies $c_T = h_T$ and the following recursive formula*

$$c_{t-1} = h_{t-1} + \lambda^2 h_{t-1}^2 c_t (1 + \phi(\xi^2) c_t). \tag{32}$$

The proof is given in Appendix E.1 in the supplementary material. Finally, setting $a_1 = 0$, we derive the formula for finding $a_{t+1}$. It follows from Eq. (30) and Eq. (31) that

$$
\begin{aligned}
a_{t+1} &= \frac{1}{2}\lambda c_{t+1} \left( \boldsymbol{x}_t - \frac{1}{2}\boldsymbol{1}_t \right)^{\mathrm{T}} \boldsymbol{z}_t \\
&= \frac{1}{2}\lambda c_{t+1} h_t \left( 2\frac{a_t}{c_t} + x_t - \frac{1}{2} \right).
\end{aligned} \tag{33}
$$

As in this section, if all $\xi_t$ are identically fixed, the approximation accuracy may be worse than the case where $\xi_t$ changes over time while the calculation time of the prediction can be reduced because it is computable by a linear form as in the last expression of $a_{t+1}$. However, if we change $\xi_t$ at every $t$, we need to calculate the inverse matrix in Eq. (24), which makes the calculation time longer. Thus, how to determine the parameter $\xi_t$ entails a trade-off between calculation time and approximation accuracy. In addition, if some prior knowledge of $\xi_t$ such as the range or the average is available, fixing all $\xi_t$ to a common value based on it enables online prediction with less calculation time while maintaining approximation accuracy. A guideline for choosing a common value $\xi$ of the parameter is discussed in Appendix C in the supplementary material.

### 3.6. Evaluation of Upper and Lower Bounds of Regret

In this subsection, we evaluate the upper and lower bounds of the regret. The analysis of the regret is intractable in the general case in Section 3.4 because we no longer have the explicit expressions of the matrix $R_T$ of Eq. (24) as in Lemma 1 nor of its recursive decomposition $R_t$. Hence, we analyze the regret for the special case in Section 3.5. Details of the derivation are shown in Appendix A. As a result, the minimax regret $R^* = \frac{1}{16}\sum_{t=1}^{T} c_t$ is upper-bounded by

$$O\left( \frac{T}{\sqrt{\lambda}} \right)$$

and lower-bounded by

$$\Omega\left( \frac{T}{\sqrt{\lambda}} \right).$$

These regret bounds hold for any variational parameter $\xi$, and are the same orders as those for the minimax online prediction under the squared loss (Koolen et al., 2015). The fact that the orders of regret bounds do not depend on $\xi$ suggests that the same regret bounds hold for the case where the variational parameter $\xi_t$ changes over time too[A1-5]. In addition, if the approximation error of the offline prediction is the order of $o(\frac{T}{\sqrt{\lambda}})$, it is conjectured that the regret of online prediction without the approximation can also be upper-bounded by $O(\frac{T}{\sqrt{\lambda}})$.

## 4. Numerical Experiments

### 4.1. Example of Online Prediction

We first show an example of online prediction for time varying Bernoulli process using the proposed method. A binary sequence of length $T = 360$ was generated from the Bernoulli process with parameter $\theta_t^* = 0.4 + 0.35 \sin(2\pi t/180)$ in trial $t$. In the proposed method, Eq. (33) is used as the prediction strategy. The parameters $\lambda = 4$ and $\xi = 0$ were used. We fixed the first predicted value $a_0 = 0$. Fig. 2 summarizes a example of offline and online predictions and the parameters of the Bernoulli process ("true_rate"). In addition, The lower part of the Fig. 2 shows the binary sequence generated according to the underlying process ("true_rate").
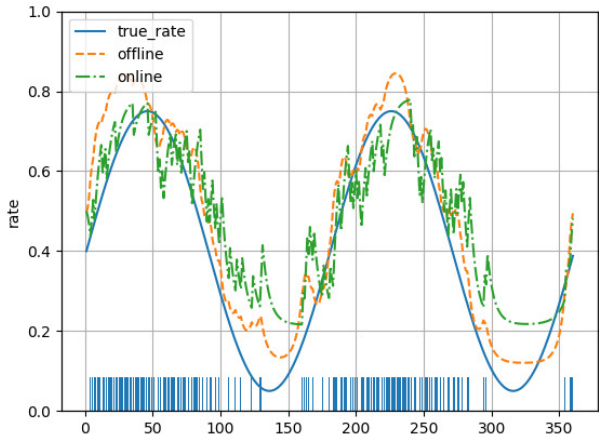


Figure 2: An example of the result of online prediction.

### 4.2. Comparison between Bayesian Prediction Method and Proposed Method

The previous study by Watanabe and Okada (2011) proposed the Bayesian prediction method which predicts the varying binomial process from the binary observation at each time point based on the Bayesian inference. The detailed procedure of the Bayesian prediction method is described in Appendix D in the supplementary material. The filtering algorithm in that method, sequentially computes the distributions of one-step ahead prediction $p(a_{t+1}|\boldsymbol{x}_t)$ and filtering $p(a_t|\boldsymbol{x}_t)$ at each time point. The mean of the one-step ahead predictive distribution can be used as the result of online prediction. Then, smoothing estimates the past states after observing the whole data sequence. The mean of the smoothing distribution can be used as the result of offline prediction. The approximation used in the Bayesian prediction method is described in detail in Section 3.1.

We compare prediction accuracies of the proposed method and Bayesian prediction method by the three measures, namely the approximate regret with Eq. (8) as the loss function $R_{\mathrm{approx}}$, the original regret with Eq. (3) as loss function $R_{\mathrm{logistic}}$ and the KL
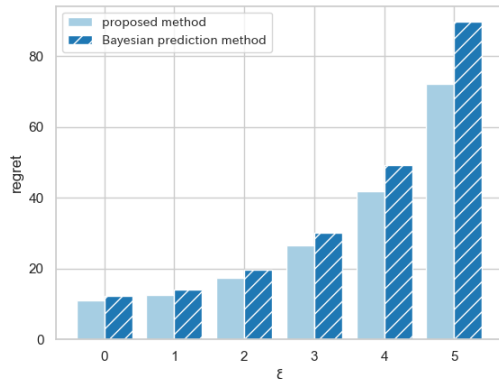
Figure 3: The maximum regret with approximation, $R_{\mathrm{approx}}$.

divergence,

$$
D_{\mathrm{KL}} = \sum_{t=1}^{T} \left\{ \theta_t^* \ln \frac{\theta_t^*}{\theta_t} + (1 - \theta_t^*) \ln \frac{1 - \theta_t^*}{1 - \theta_t} \right\},
$$

where $\theta_t^*$ is the true underlying probability and $\theta_t = a^{-1}(a_t)$ is the predicted probability. The offline prediction problem in Eq. (10) corresponds to the maximum a posteriori estimation in the Bayesian method (Watanabe and Okada, 2011). We used the one-step ahead prediction for online prediction (Watanabe and Okada, 2011). The proposed method uses Eq. (33) for online prediction. 500 binary sequences of length $T = 360$ were generated from the Bernoulli process with parameter $\theta_t^* = 0.35 + 0.3 \sin(2\pi t/180)$ in trial $t$. We fixed the variational parameter $\xi$ to 0, 1, 2, 3, 4 and 5 in both methods, predicted each of 500 binary sequences to compute $R_{\mathrm{approx}}$, $R_{\mathrm{logistic}}$ and $D_{\mathrm{KL}}$ and compared the maximum value of each measures between the methods. In addition, the parameter $\lambda$ was obtained for each $\xi$ by the average of the values estimated for each binary sequences in the process of offline prediction of Bayesian prediction, and we let the first predicted value $a_0 = 0$. Note that $\xi$ is set more advantageously for Bayesian prediction than for the proposed method. We show in Fig. 3 ($R_{\mathrm{approx}}$), Fig. 4 ($R_{\mathrm{logistic}}$) and Fig. 5 ($D_{\mathrm{KL}}$) the maximum values of the respective measures.

Fig. 3 demonstrates that the maximum regret is smaller in the proposed method than in the Bayesian prediction method regardless of the parameter $\xi$. Even though this is a result to a limited number of input sequences, this represents the validity of the guarantee that the proposed method minimizes the maximum regret under approximation. Also, Fig. 4 shows that the maximum regret with the original logistic loss function are comparable for $\xi = 0$ and 1. Since this is the maximum regret without approximation, the guarantee of the proposed method is not effective rigorously. For $\xi \geq 2$, however, the proposed method outperforms the Bayesian method in terms of the maximum original regret (Fig. 4) and the KL divergence from the true process (Fig. 5). This implies that if the approximation is made by an appropriate $\xi$ (even if it is not optimal), the proposed method can reduce the maximum regret more than the Bayesian prediction method. In addition, we can see a
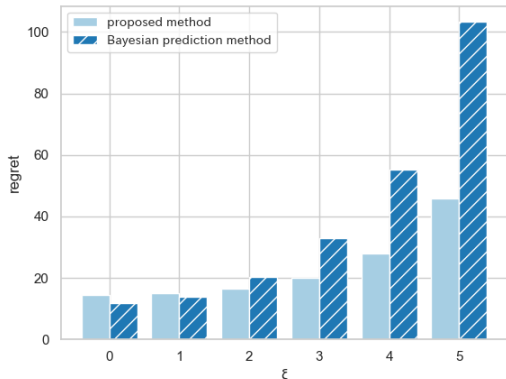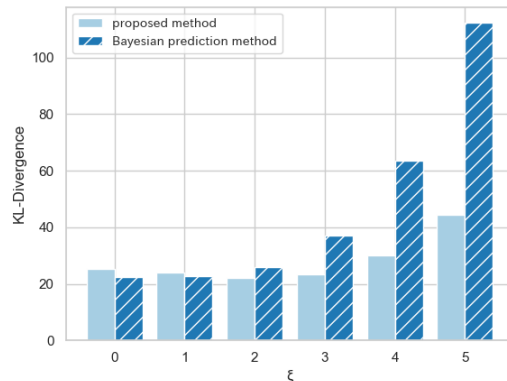
Figure 4: The maximum regret without approximation, $R_{\text{logistic}}$.

Figure 5: The maximum KL divergence, $D_{\text{KL}}$.

similar tendency in Fig. 4 and Fig. 5, which suggests that reducing the maximum regret leads to the improvement in the difference between the true variation and prediction. It is desirable to determine the variational parameter $\xi$ appropriately in the proposed method as well as in the Bayesian prediction method. If we have some prior knowledge on the true underlying process such as the average and the maximum variation (from probability 0.5), we can use it to determine $\xi$ because $\xi$ has a one-to-one correspondence with the value of Bernoulli probability (Jaakkola and Jordan, 2000). The influence of the parameter $\xi$ on prediction is discussed in Appendix B in the supplementary material. As discussed in the last part of Appendix C in the supplementary material, setting $\xi$ to the value corresponding to the maximum variation from probability 0.5 (corresponding to $\xi = 0$) may be a guideline for choosing $\xi$. The corresponding value for the case of this experiment is $\xi = 2.944$, for which we have seen that the proposed method is likely to outperform the Bayesian method in terms of the maximum regret (see the case of $\xi = 3$ in Figs. 3–5).

## 5. Conclusion

In this paper, we developed an online prediction method of varying Bernoulli process from a single binary sequence by using variational approximation for the logistic error function, which has a guarantee to minimize the maximum regret under approximation. The upper and lower bounds of the minimax regret were evaluated theoretically. We discussed the influence of the approximation on the regret and the difference between the true varying Bernoulli process and the prediction through numerical experiments. Although the proposed method has only a guarantee about the regret under approximation, it is suggested that the influence of the approximation on the regret can be reduced by properly setting the variational parameter. Our future directions include devising methods to determine the variational parameter and the regularization parameter for smoothness, and the theoretical verification of the influence by the variational approximation.

## Acknowledgments

## References

Junyu Cao, Wei Sun, and Zuo-Jun Max Shen. Sequential choice bandits: Learning with marketing fatigue. *SSRN Electronic Journal*, January 2019.

Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games.* Cambridge University Press, 2006.

John P. Cunningham, Vikash Gilja, Stephen I. Ryu, and Krishna V. Shenoy. Methods for estimating neural firing rates, and their application to brain-machine interfaces. *Neural Networks*, 22(9):1235 – 1246, 2009.

Te Sun Han and Kingo Kobayashi. *Mathematics of information and coding.* American Mathematical Society, 2007.

Mark Herbster and Manfred K. Warmuth. Tracking the best linear predictor. *Journal of Machine Learning Research*, 1:281–309, 2001.

GY Hu and Robert F. O'Connell. Analytical inversion of symmetric tridiagonal matrices. *Journal of Physics A: Mathematical and General*, 29(7):1511, 1996.

Tommi S. Jaakkola and Michael I. Jordan. Bayesian parameter estimation via variational methods. *Statistics and Computing*, 10(1):25–37, January 2000.

Rudolf Kalman. A new approach to linear filtering and prediction problems. *Transactions of the ASME - Journal of basic Engineering*, 82:35–45, 01 1960.

Wouter M. Koolen, Alan Malek, and Peter L. Bartlett. Efficient minimax strategies for square loss games. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3230–3238. Curran Associates, Inc., 2014.

Wouter M. Koolen, Alan Malek, Peter L. Bartlett, and Yasin Abbasi. Minimax time series prediction. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2557–2565. Curran Associates, Inc., 2015.

Shinsuke Koyama and Shigeru Shinomoto. Empirical Bayes interpretations of random point events. *Journal of Physics A: Mathematical and General*, 38(29):L531–L537, jul 2005.

Shuai Li, Claudio Gentile, Alexandros Karatzoglou, and Giovanni Zappella. Data-dependent clustering in exploration-exploitation algorithms. *Computing Research Repository*, abs/1502.03473, 2015. URL http://arxiv.org/abs/1502.03473.

Edward Moroshko and Koby Crammer. Weighted last-step min-max algorithm with improved sub-logarithmic regret. *Theoretical Computer Science*, 558(13):107–124, 2014.

Jorma Rissanen and Glen Langdon. Universal modeling and coding. *IEEE Transactions on Information Theory*, 27(1):12–23, January 1981.

Ken Takiyama and Masato Okada. Switching state space model for simultaneously estimating state transitions and nonstationary firing rates. In *Advances in Neural Information Processing Systems 23*, pages 2271–2279. Curran Associates Inc., 2010.

Kazuho Watanabe and Masato Okada. Approximate Bayesian estimation of varying binomial process. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, E94.A(12):2879–2885, 2011.

## Appendix A. Upper and Lower Bounds of Minimax Regret

### A.1. Upper Bound

To bound the minimax regret $R^*$ from above, we evaluate $c_t$ from above. It follows from Eq. (32) that
$$c_{t-1} = h_{t-1} + \lambda^2 h_{t-1}^2 c_t (1 + \phi(\xi^2)c_t).$$
In addition, since $h_t \le \frac{h}{\phi(\xi^2)}$, if we put

$$c'_{t-1} = \frac{h}{\phi(\xi^2)} + \lambda^2 \frac{h^2}{\phi(\xi^2)^2} c'_t (1 + \phi(\xi^2)c'_t),$$

$c'_t \ge c_t$ holds. While $c'_t$ grows as $t$ decreases, it has its limit. Let $C$ be the limit of $c'_t$. It satisfies the following equation:

$$C = \frac{h}{\phi(\xi^2)} + \lambda^2 \frac{h^2}{\phi(\xi^2)^2} C (1 + \phi(\xi^2)C). \tag{34}$$

we solve Eq. (34) for $C$. Putting $\frac{\lambda}{\phi(\xi^2)} = \tilde{\lambda}$, we have

$$C = \frac{h}{\phi(\xi^2)} + \tilde{\lambda}^2 h^2 C + \tilde{\lambda}^2 h^2 \phi(\xi^2)C^2$$

$$\Leftrightarrow 0 = \tilde{\lambda}^2 h^2 \phi(\xi^2)C^2 + (\tilde{\lambda}^2 h^2 - 1)C + \frac{h}{\phi(\xi^2)}$$

$$\Leftrightarrow C = \frac{1 - \tilde{\lambda}^2 h^2 \pm \sqrt{(\tilde{\lambda}^2 h^2 - 1)^2 - 4\tilde{\lambda}^2 h^3}}{2\phi(\xi^2)\tilde{\lambda}^2 h^2}$$

In the two solutions of $C$, when we bound $c_t$ from above, the smaller $C$ provides tighter, bound

$$C = \frac{\frac{1}{h^2} - \tilde{\lambda}^2 - \sqrt{\left(\tilde{\lambda}^2 - \frac{1}{h^2}\right)^2 - 4\tilde{\lambda}^2 \frac{1}{h}}}{2\phi(\xi^2)\tilde{\lambda}^2}. \tag{35}$$

Substituting Eq. (29) for Eq. (35), we have

$$C = \frac{\frac{1 + 4\tilde{\lambda} + \sqrt{1+4\tilde{\lambda}} + 2\tilde{\lambda}\sqrt{1+4\tilde{\lambda}}}{2}}{2\phi(\xi^2)\tilde{\lambda}^2} - \frac{\sqrt{\frac{1 + 8\tilde{\lambda} + 12\tilde{\lambda}^2 + \sqrt{1+4\tilde{\lambda}} + 6\tilde{\lambda}\sqrt{1+4\tilde{\lambda}} + 4\tilde{\lambda}^2\sqrt{1+4\tilde{\lambda}}}{2}}}{2\phi(\xi^2)\tilde{\lambda}^2} \tag{36}$$

We evaluate the order on $\tilde{\lambda}$ for each term of the denominator and numerator in Eq. (36). The order of the denominator is $O(\tilde{\lambda}^2)$, the first term of the numerator is $O\left(\tilde{\lambda}\sqrt{\tilde{\lambda}}\right) = O(\tilde{\lambda}^{\frac{3}{2}})$, all the second term of the numerator is $\left(\sqrt{\tilde{\lambda}^2\sqrt{\tilde{\lambda}}}\right) = O(\tilde{\lambda}^{\frac{5}{4}})$. Therefore, the order of $C$ is

$$O\left(\frac{\tilde{\lambda}^{\frac{3}{2}}}{\tilde{\lambda}^2}\right) = O\left(\frac{1}{\sqrt{\tilde{\lambda}}}\right).$$

In addition, since $\xi$ is a constant, the regret is upper-bounded by $\frac{1}{16}\sum_{t=1}^{T} c_t = O\left(\frac{T}{\sqrt{\lambda}}\right)$.

### A.2. Lower Bound

To bound the minimax regret $R^*$ from below, we evaluate $c_t$ from below. It follows from Eq. (32) that

$$c_{t-1} = h_{t-1} + \lambda^2 h_{t-1}^2 c_t (1 + \phi(\xi^2) c_t).$$

Since $\lambda, h_t$ and $c_t$ are all non-negative, ignoring the square of $c_t$ yields

$$c_{t-1} \geq h_{t-1} + \lambda^2 h_{t-1}^2 c_t.$$

Let $r_t = 1 - (\lambda h)^{2t}$,

$$h_t = \frac{h r_t}{r_{t+1}} \tag{37}$$

holds. Considering recursion and using Eq. (37), we can bound $c_t$ from below,

$$c_t \geq h \sum_{k=t}^{T} (\lambda h)^{2(k-t)} \frac{r_t^2}{r_k r_{k+1}},$$

as detailed in Appendix E.4 in the supplementary material. Since $\frac{r_t^2}{r_i r_{i+1}}$ is a decreasing function on $i$, $\frac{r_t^2}{r_i r_{i+1}} \geq \frac{r_t}{r_{t+1}}$ holds and it follows that

$$\sum_{t=1}^{T} c_t \geq h \sum_{t=1}^{T}\sum_{k=t}^{T} (\lambda h)^{2(k-t)} \frac{r_t}{r_{t+1}} \geq h \int_0^{T-1}\int_{t+1}^{T} (\lambda h)^{2(k-t)} \frac{r_t}{r_{t+1}} dk dt = \Omega\left(-\frac{hT}{2\log(\lambda h)}\right).$$

The derivation of the last step is detailed in Appendix E.5 in the supplementary material. Here, since

$$-\ln(\lambda h) = -\ln\left(\frac{2}{\frac{1}{\lambda} + \frac{2}{\phi(\xi^2)} + \frac{1}{\lambda}\sqrt{1 + 4\frac{\lambda}{\phi(\xi^2)}}}\right) = \Omega\left(\frac{1}{\sqrt{\lambda}}\right),$$

$$h = \frac{2}{1 + \frac{2\lambda}{\phi(\xi^2)} + \sqrt{1 + 4\frac{\lambda}{\phi(\xi^2)}}} = \Omega\left(\frac{1}{\lambda}\right),$$

the regret is lower-bounded by $\frac{1}{16}\sum_{t=1}^{T} c_t = \Omega\left(\frac{T}{\sqrt{\lambda}}\right)$.