

VARIABLE FORWARD REGULARIZATION TO REPLACE RIDGE IN ONLINE LINEAR REGRESSION

Anonymous authors

Paper under double-blind review

ABSTRACT

Forward regularization (-F) [Azoury & Warmuth \(2001\)](#) with unsupervised knowledge was proposed to replace canonical *Ridge regularization* (-R) in online linear learners, which achieves lower relative regret bounds ([Della Vecchia & Basu, 2023](#)). However, we observe that -F cannot perform as expected in practice, even possibly losing to -R for online learning tasks. We identify two main causes for this: (1) inappropriate intervention penalty; (2) potential non-i.i.d nature in online learning, both of which result in unstable posterior distribution and optimal offset of the learner. To improve these, we propose *Variable Forward regularization* (- k F), a more general style with -F intensity modulated by a variable k . We further derive - k F algorithm to online learning tasks, which shows holistic recursive closed-form updates and superior performance compared to both -R and -F. Moreover, we theoretically establish the relative regrets of - k F in online learning, showing that it has a tighter upper bound than -F in adversarial settings. We also introduce an adaptive - k F, termed - k F-Bayes, to curb unstable penalties caused by non-i.i.d and mitigate intractable tuning of hard k based on Bayesian learning for online learning. In experiments, we adapted - k F and - k F-Bayes into class incremental scenario, where it realized less forgetting and non-replay. Results distinctly demonstrate the efficacy of using - k F and - k F-Bayes.

1 INTRODUCTION

Unlike offline training, the online learning (OL) scenario emphasizes progressive model updates and delivers immediate decision-making on non-stationary data/task streams. Applications of OL methodology include chatbots ([Chowdhury et al., 2023](#)), injection attack detection ([Toyer et al., 2024](#)), and recommendation systems ([Jiang et al., 2024](#)) of scenes with real-time demands. During the OL process, the backbone is expected to preserve previous knowledge without past replay and reduce learning regrets ([Buening et al.](#)). In this paper, we study linear regression, one basic element of connectionist models, equipped with a novel regularization and its characteristics in online continual learning (CL) scenarios.

Forward regularization (-F) proposed in ([Azoury & Warmuth, 2001](#); [Vovk, 2001](#)) was proven to have better regret bounds compared to using canonical *Ridge regularization* (-R) in linear regression with challenging adversarial bounded observations during OL [Della Vecchia & Basu \(2023\)](#). Specifically, -F is realized by adding an extra Frobenius-norm penalty containing upcoming unsupervised knowledge to -R, which leads the gradients of learnable weights to be more concerned with future prediction. Besides the decrease in regret, the transductive method -F has low computation complexity and regulates the learning rate in OL. As a result, -F is of practical relevance as it enhances prediction accuracy by integrating unlabeled data.

However, we observe that -F could not perform as expected in experiments, even possibly failing to -R during OL. We argue the reason is two-fold. (1) *Improper regularization*. -F defaults the penalty intensity of the unsupervised term is constant 1.0, resulting in improper penalty intervention in special tasks. (2) *Non-i.i.d disturbance*. Premise i.i.d is not held in OL, which incurs unstable regularization due to varying distribution and batch volume. This causes residuals of empirical risk minimization (ERM) and model optimum offsets, especially in learning long task streams.

To solve the problems, we first propose the *Variable Forward regularization* (- k F), adding a hard threshold k to control the penalty strength for different tasks. The more general - k F is theoretically

054 proven to have better performance and tighter regret bounds than -R and -F by two ways if the k can
 055 be given properly. The $-kF$ also generates one-shot closed-form incremental updates and variable
 056 learning rate, and has less learning dissipation. Moreover, to curb unstable penalties caused by non-
 057 i.i.d and mitigate intractable tuning of hard k during OL, we improve the $-kF$ to $-kF$ -Bayes style,
 058 further enabling the soft k to be adaptively determined based on Bayesian learning and estimated
 059 distribution.

060 Our algorithm is assessed within the framework of Continual Learning (CL), which is a sub-
 061 set of OL. CL involves learning a sequence of Q distinct tasks or episodes, denoted as $\mathcal{T} =$
 062 $\{\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_Q\} | \mathcal{T}_q = (\mathcal{X}_q, \mathcal{Y}_q)\}$, and includes unseen classes (French, 1999). Among the various
 063 scenarios of CL, Class Incremental Learning (CIL) poses the greatest challenge. In CIL, the learner
 064 is required to integrate and retain knowledge of all previous classes, represented as $\mathcal{Y} = \bigcup_{q=1}^Q \mathcal{Y}_q$
 065 without access to task identities during testing. To overcome catastrophic forgetting, EWC Kirk-
 066 patrick et al. (2017) is a pioneer regularization-based work, restricting previous weights from wrong
 067 skewing by imposing penalties. Many new methods, such as Online EWC Schwarz et al. (2018),
 068 SI Zenke et al. (2017), CRNet Li & Zeng (2023), RanPAC McDonnell et al. (2024), NICE Gurbuz
 069 et al. (2024) are born recent years. However, CIL assumes task-wise i.i.d and task boundary hinting
 070 update opportunity, hardly meeting practical needs. The harsher online task-free CIL (OTCIL) is
 071 expected to generate immediate decision-making with neither i.i.d nor boundary index in each task,
 072 where one task comes in non-stationary batches rather than all at once in CIL. Related works include
 073 GEM Lopez-Paz & Ranzato (2017), GSS Aljundi et al. (2019), and DYSON (He et al., 2024). We
 074 also applied the proposed $-kF$ and $-kF$ -Bayes methodologies to OTCIL scenarios and provided an
 075 in-depth analysis.

076 Contributions are summarized as follows:

- 077 1. **Propose $-kF$ and derive regret bound:** To control the intervention rate of unsupervised
 078 knowledge for improved performance, -F is extended to the more general $-kF$, incorporat-
 079 ing a k -factored penalty. This enhancement aims to refine the learning gradients in linear
 080 regression and reduce regrets in OL. We formulate the $-kF$ algorithm with variable learning
 081 rates and recursive closed-form updates, and demonstrate that it prevents learning dissipation
 082 while achieving tighter relative regret bounds compared to -F and -R.
- 083 2. **Enhance $-kF$ with Bayes:** To curb unstable penalties caused by non-i.i.d and mitigate
 084 intractable tuning of k of $-kF$ during OL, we further propose $-kF$ -Bayes, enabling the soft
 085 k to be adaptively determined based on Bayesian learning and distribution estimation.
- 086 3. **Practical application:** We integrated $-kF$ and $-kF$ -Bayes into randomized learners and
 087 evaluated them in (OT)CIL scenarios using both tabular and image datasets. The results
 088 demonstrated the effectiveness of our methods.

090 2 PRELIMINARY AND ADVERSARIAL REGRET BOUNDS

092 **Lemma 1. Bregman divergence** is used to measure relative projection distance between distribu-
 093 tions (Azoury & Warmuth, 2001). For a real-valued differentiable convex projection $G : \Theta \rightarrow$
 094 \mathbb{R} , Bregman divergence Δ is defined as:

$$095 \Delta_G(\tilde{\theta}, \theta) := G(\tilde{\theta}) - G(\theta) - (\tilde{\theta} - \theta)^T \nabla_{\theta} G(\theta), \quad (1)$$

097 where the θ is a vector, and ∇_{θ} denotes the gradient operator on vector θ .

098 **Lemma 2. Bregman divergence property.**

- 099 1. The divergence is a linear operator. $\forall \mu \geq 0, \Delta_{G_1 + \mu G_2}(\tilde{\theta}, \theta) = \Delta_{G_1}(\tilde{\theta}, \theta) + \mu \Delta_{G_2}(\tilde{\theta}, \theta)$.
- 100 2. If $G_1(\theta) - G_2(\theta) = \omega^T \theta + v, \omega \in \mathbb{R}^{|\Theta|}, v \in \mathbb{R}$, then $\Delta_{G_1}(\tilde{\theta}, \theta) = \Delta_{G_2}(\tilde{\theta}, \theta)$.

102 **Setup:** following the classical setups of linear regression in OL Azoury & Warmuth (2001), a
 103 representation steam $\mathcal{X} = \{\mathbf{x}_t\}_{t=1}^T \subseteq \mathbb{R}^d$ paired with corresponding observations $\mathcal{Y} = \{y_t\}_{t=1}^T \subseteq$
 104 \mathbb{R} is received by the learner f one-by-one in trials $t = 1, 2, \dots, T$. Joint domain distribution $\mathcal{X} \times \mathcal{Y} \sim \mathcal{P}$
 105 but drifting occurs across batches. In -R, f_t is updated on visible (\mathbf{x}_t, y_t) and evolves to f_{t+1} , hoping
 106 to predict \hat{y}_{t+1} close to y_{t+1} based on \mathbf{x}_{t+1} . The oracle learner is defined as:

$$107 f : y_t = \mathbf{x}_t^T \theta_t^* + \varepsilon_t \quad \exists \theta_t^* \in \Theta \subseteq \mathbb{R}^d, 1 \leq t \leq T \quad (2)$$

where the ε_t is Gaussian noise, θ^* can be learnable weights, and Θ is weight space.

In the optimization target, the initial G , incurred loss on x_t , and forward predictive loss are respectively designated to $U_0(\theta) = \frac{1}{2}\theta^T\eta_0^{-1}\theta$, $\mathcal{L}_t(\theta) = \frac{1}{2}\|x_t^T\theta - y_t\|_2^2$, and $\hat{\mathcal{L}}_{t+1}(\theta) = \frac{1}{2}\|x_{t+1}^T(\theta - \theta_0)\|_2^2$, where η_0^{-1} is a symmetric positive definite matrix and θ_0 is the initial parameter.

Lemma 3. Offline learning refers to the learning behavior of an expert on global tasks. Online learner is contrasted to the expert by relative regret bounds. Assume solutions always exist in $\theta \in \Theta$:

$$\theta_{Q+1} = \operatorname{argmin}_{\theta} U_{Q+1}(\theta), \quad (3)$$

where $U_{Q+1}(\theta) = \Delta_{U_0}(\theta, \theta_0) + \mathcal{L}_{1..Q}(\theta)$, θ_{Q+1} represents the finally updated parameter for future predictions after the last Q -th task knowledge acquisition completed.

Lemma 4. Online-to-offline regret bounds are defined as the upper bounds of cumulative regrets of an online learner over those of the offline expert, which quantify the gap to the best expert and be regarded as the cost of hiding future data from the learner. The upper bound of the learner using -R:

$$\sum_{t=1}^T \mathcal{L}_t(\theta_t^r) - \min_{\theta} (\frac{1}{2}\lambda\|\theta - \theta_0\|_2^2 + \sum_{t=1}^T \mathcal{L}_t(\theta)) \leq 2Y_m^2 \operatorname{dIn}(\frac{TX_m^2}{\lambda} + 1) \quad (4)$$

where $X_m = \max_{1 \leq t \leq T} \{\|x_t\|_{\infty}\}$, $Y_m = \max_{1 \leq t \leq T} \{|y_t|, |x_t^T \theta_t|\}$, and λ is the penalty factor. This is based on adversarial setups and highlights the prediction $x_t^T \theta_t \in [-Y_m, Y_m]$ for history value restraints.

The upper bound of the learner using -F:

$$\sum_{t=1}^T \mathcal{L}_t(\theta_t^f) - \min_{\theta} (\frac{1}{2}\lambda\|\theta\|_2^2 + \sum_{t=1}^T \mathcal{L}_t(\theta)) \leq \frac{1}{2}Y_m^2 \operatorname{dIn}(\frac{TX_m^2}{\lambda} + 1) \quad (5)$$

where $X_m = \max_{1 \leq t \leq T} \{\|x_t\|_{\infty}\}$, $Y_m = \max_{1 \leq t \leq T} \{|y_t|\}$. Prediction range assumption is removed and only retains $y_t \in [-Y_m, Y_m]$. Note the -F bound is at least 4 times better than -R's.

3 PROPOSED METHODOLOGY AND REGRET BOUND DERIVATION

3.1 ONLINE LEARNING FRAMEWORK FOR $-kF$ STYLE

Theorem 1. Incremental offline learning using $-kF$. We define the $-kF$ optimization target and express an incremental offline pattern here. For $0 \leq t \leq T$, the following equations are optimized:

$$\theta_{t+1} = \operatorname{argmin}_{\theta} U_{t+1}(\theta) \quad (6)$$

$$U_{t+1}(\theta) = \Delta_{U_0}(\theta, \theta_0) + \mathcal{L}_{1..t}(\theta) + k \cdot \hat{\mathcal{L}}_{t+1}(\theta)$$

where $\hat{\mathcal{L}}$ is the estimated loss on upcoming data (reserved prior knowledge is also possible). $-kF$ integrates unsupervised knowledge and affects optimization gradients (also a posterior optimum) when updating the present model. Note $\theta_1 = \operatorname{argmin}_{\theta} U_1(\theta) = \theta_0$ when $t = 0$.

Theorem 2. OL using $-kF$. For $0 \leq t \leq T$, with previous $\mathcal{L}_{1..t-1}$ concentrated into Bregman divergence ΔU_t , we define the OL process of $-kF$ as optimizing a recursive equation without replay.

$$\theta_{t+1} = \operatorname{argmin}_{\theta} \Delta U_t(\theta, \theta_t) + \mathcal{L}_t(\theta) + k \cdot \hat{\mathcal{L}}_{t+1}(\theta) - k \cdot \hat{\mathcal{L}}_t(\theta) \quad (7)$$

When $0 < k \leq 1$, it can be regarded as partial forward knowledge intervention because ERM on x_t is larger than the forward penalty of x_{t+1} . Having $k > 1$ leads to an excess in forward regularization.

Proof: Please see Appendix A.1

3.2 ONLINE LEARNING ALGORITHM FOR $-kF$ STYLE

Theorem 3. OL algorithm using $-kF$. The step-wise updates of learnable weight θ in $-kF$ follow:

$$\theta_{t+1} = \theta_t - \eta_{t+1} [(x_t x_t^T + k \cdot x_{t+1} x_{t+1}^T - k \cdot x_t x_t^T) \theta_t - x_t y_t] \quad (8)$$

where k is a constant, $\theta_0 = 0$, symmetric positive defined $\eta_0 = (\lambda \cdot \mathbf{I})^{-1}$, and variable learning rates $\eta_{t+1} = (\eta_0^{-1} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^T + k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T)^{-1}$. The step-wise updates of η_{t+1} follow:

$$\begin{aligned} \eta_{t+1}^\dagger &= \eta_t^\dagger - \eta_t^\dagger \mathbf{x}_t (\mathbf{I} + \mathbf{x}_t^T \eta_t^\dagger \mathbf{x}_t)^{-1} \mathbf{x}_t^T \eta_t^\dagger \\ \eta_{t+1} &= \eta_{t+1}^\dagger - \eta_{t+1}^\dagger k \cdot \mathbf{x}_{t+1} (\mathbf{I} + k \cdot \mathbf{x}_{t+1}^T \eta_{t+1}^\dagger \mathbf{x}_{t+1})^{-1} \mathbf{x}_{t+1}^T \eta_{t+1}^\dagger \end{aligned} \quad (9)$$

where $\eta_{t+1}^\dagger = (\eta_0^{-1} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^T)^{-1}$ is maintained as an intermediate variable for computation. This OL process of using $-k\text{F}$ is non-replay, and the Bregman divergence term changes over time.

Proof: Please see Appendix A.2

Remark 1. Our $-k\text{F}$ style is more general. If $k = 0$, Theorem 2 degenerates to -R style:

$$\begin{aligned} \theta_{t+1} &= \operatorname{argmin}_{\theta} \Delta_{U_t}(\theta, \theta_t) + \mathcal{L}_t(\theta) \\ \theta_{t+1} &= \theta_t - \eta_{t+1} [\mathbf{x}_t \mathbf{x}_t^T \theta_t - \mathbf{x}_t y_t] \end{aligned} \quad (10)$$

If $k = 1$, Theorem 2 degenerates to -F style:

$$\begin{aligned} \theta_{t+1} &= \operatorname{argmin}_{\theta} \Delta_{U_t}(\theta, \theta_t) + \mathcal{L}_t(\theta) + \hat{\mathcal{L}}_{t+1}(\theta) - \hat{\mathcal{L}}_t(\theta) \\ \theta_{t+1} &= \theta_t - \eta_{t+1} [\mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \theta_t - \mathbf{x}_t y_t] \end{aligned} \quad (11)$$

This further proves $-k\text{F}$ style is more general and the performance of using $-k\text{F}$ must not be inferior to that of -R or -F if k is properly designed, and such k values exist with a high probability.

Remark 2. We further derive several equalities before the regret bound derivation. Given $\eta_{t+1}^{-1} = \eta_t^{-1} + (1 - k) \cdot \mathbf{x}_t \mathbf{x}_t^T + k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T$ and $\theta_0 = 0$, for $1 \leq t \leq T$ we have:

$$\theta_{t+1} = \eta_{t+1} (\eta_t^{-1} \theta_t + \mathbf{x}_t y_t) \quad (12)$$

$$\theta_{t+1} = \theta_t - \eta_t [(\mathbf{x}_t \mathbf{x}_t^T + k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T - k \cdot \mathbf{x}_t \mathbf{x}_t^T) \theta_{t+1} - \mathbf{x}_t y_t] \quad (13)$$

$$\theta^{off} = (\eta_0^{-1} + \sum_{i=1}^T \mathbf{x}_i \mathbf{x}_i^T)^{-1} \eta_{T+1}^{-1} \theta_{T+1} \quad (14)$$

Proof: Please see Appendix A.3

3.3 RELATIVE REGRET BOUND FOR $-k\text{F}$ STYLE

Theorem 4. Relative Regrets for $-k\text{F}$. For any $\mathcal{X} = \{\mathbf{x}_t\}_{t=1}^T$ and any $\theta \in \Theta$, the cumulative regrets of the online learner using $-k\text{F}$ style relative to the offline expert can be calculated as follows:

$$\begin{aligned} & \sum_{t=1}^T \mathcal{L}_t(\theta_t) - \min_{\theta \in \theta^{off}} (\Delta_{U_0}(\theta, \theta_0) + \sum_{t=1}^T \mathcal{L}_t(\theta)) \\ &= \sum_{t=1}^T [\Delta_{U_{t+1}}(\theta_t, \theta_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\theta_t) + k \cdot \hat{\mathcal{L}}_t(\theta_t)] + k \cdot \hat{\mathcal{L}}_{T+1}(\theta) - \Delta_{U_{T+1}}(\theta, \theta_{T+1}) \end{aligned} \quad (15)$$

Proof: Please see Appendix A.4

Theorem 5. Relative Regret Bounds for $-k\text{F}$. The upper bound of the learner using $-k\text{F}$:

$$\mathbb{E}[\sum_{t=1}^T \mathcal{L}_t(\theta_t) - \min_{\theta \in \theta^{off}} (\Delta_{U_0}(\theta, \theta_0) + \sum_{t=1}^T \mathcal{L}_t(\theta))] = \frac{k}{2} Y_m^2 d \ln(1 + \frac{TX_m^2}{\lambda + (k-1) \cdot X_m^2}) \quad (16)$$

where $X_m = \max_{1 \leq t \leq T} \{\|\mathbf{x}_t\|_\infty\}$, $Y_m = \max_{1 \leq t \leq T} \{y_t\}$, and $k > 0$. Like (5), (16) still maintains the advantages of border (i.e. $\mathbf{x}_t^T \theta_t \in [-Y_m, Y_m]$) removal. (16) corresponds to (5)'s form if $k = 1$.

Proof: Please see Appendix A.5

Remark 3. Let $\mathbb{E}[\sum_{t=1}^T \mathcal{L}_t(\theta_t) - \min_{\theta \in \theta^{off}} (\Delta_{U_0}(\theta, \theta_0) + \sum_{t=1}^T \mathcal{L}_t(\theta))] = \mathbb{E}[\Gamma]$. Assume the entire \mathcal{X} respects independent identical Gaussian distribution in proof, which allows us to involve approximation $\mathbb{E}[\mathbf{x}_{t+1}^T \eta_t \mathbf{x}_t] = 0$. Such that the Γ is a complex random variable associated with data distribution $\mathcal{P}_{\mathcal{X}}$. Consequently, the regret bounds are fluctuating. How to choose the k is a game problem, as shown in (16). One can explore the Γ distribution and lead to the proper k 's range based on a specific data environment.

Remark 4. The actual $\mathbb{E}[\Gamma]$ can be lower than that in (16), as we omit some negative penalty terms in the calculation of (55). It will be complicated to obtain the proper k 's range by studying $\frac{1}{2}Y_m^2 d\ln(\frac{TX_m^2}{\lambda} + 1) - \frac{k}{2}Y_m^2 d\ln(1 + \frac{TX_m^2}{\lambda + (k-1) \cdot \hat{X}_m^2}) \geq 0$. To avoid this, we study the growth rates of cumulative relative regrets:

$$\frac{\partial \frac{1}{2} \hat{Y}_m^2 d\ln(\frac{t\hat{X}_m^2}{\lambda} + 1)}{\partial t} = \frac{1}{2} \hat{Y}_m^2 d \cdot \frac{\hat{X}_m^2}{\lambda + t\hat{X}_m^2} \quad (17)$$

$$\frac{\partial \frac{k}{2} \hat{Y}_m^2 d\ln(1 + \frac{t\hat{X}_m^2}{\lambda + (k-1) \cdot \hat{X}_m^2})}{\partial t} = \frac{k}{2} \hat{Y}_m^2 d \cdot \frac{\hat{X}_m^2}{\lambda + (k-1) \cdot \hat{X}_m^2 + t\hat{X}_m^2} \quad (18)$$

Let (17)-(18) > 0, for $1 \leq t \leq T$, we have:

$$0 < k < 1 \quad (19)$$

where $\hat{X}_m = \max_{1 \leq i \leq t+1} \{\|\mathbf{x}_i\|_\infty\}$, $\hat{Y}_m = \max_{1 \leq i \leq t} \{|y_i|\}$.

In conclusion, the k that makes $-kF$ have a tighter regret bound than $-F$'s exists.

3.4 ADAPTIVE $-kF$ -Bayes STYLE

Potential non-i.i.d nature. In Lemma 3 and Theorem 1, impacts due to potential non-i.i.d are negligible as all training data is available, and the severity of penalties entirely maintains balance at time t . Inversely, the effects of non-i.i.d should be taken into account in Theorem 2 and 3, because we emphasize OL and CL contexts including non-i.i.d data batch, off-diagonal covariance matrix, invisible future data, and forbidden past replay. Merely setting sensitive k as fixed results in penalty imbalance and offsets the learner's optima during OL/CL process.

$-kF$ uses hard threshold k . The k mediating contributed degrees of $k \cdot \hat{\mathcal{L}}_{t+1}(\theta) - k \cdot \hat{\mathcal{L}}_t(\theta)$ determines performance of the learner using $-kF$ style. Using gradient descent or evolutionary algorithm to adjust k at each step is unreasonable due to time-consuming and non-replay, and a compromise hard k instead of synchronized one can be thrown even if global optimization for k is allowed.

Although we present a range of k in (19), it still has a relatively loose upper limitation with respect to k because we use \hat{X}_m in adversarial environment and approximation in Theorem 5, which cannot generate specific values or accurate range. To improve this, we regard the Theorem 2 and Theorem 3 as Bayesian learning processes, and study how to suppress non-i.i.d impact on posterior distribution by adaptively determined k , based on distribution estimation and penalty balance of prior terms.

Theorem 6. Adaptive $-kF$ -Bayes style The step-wise updates of learnable weight θ in $-kF$ -Bayes follow:

$$\theta_{t+1} = \theta_t - \eta_{t+1}[(\mathbf{x}_t \mathbf{x}_t^T + k_{t+1} \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T - k_t \cdot \mathbf{x}_t \mathbf{x}_t^T) \theta_t - \mathbf{x}_t y_t] \quad (20)$$

where k is a variable, $\theta_0 = 0$, symmetric positive defined $\eta_0 = (\lambda \cdot \mathbf{I})^{-1}$, and variable learning rates $\eta_{t+1} = (\eta_0^{-1} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^T + k_{t+1} \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T)^{-1}$. The step-wise updates of η_{t+1} follow:

$$\begin{aligned} \eta_{t+1}^\dagger &= \eta_t^\dagger - \eta_t^\dagger \mathbf{x}_t (\mathbf{I} + \mathbf{x}_t^T \eta_t^\dagger \mathbf{x}_t)^{-1} \mathbf{x}_t^T \eta_t^\dagger \\ k_{t+1} &= k_t = \mathbf{x}_t^T \eta_t \mathbf{x}_t \end{aligned} \quad (21)$$

$$\eta_{t+1} = \eta_{t+1}^\dagger - \eta_{t+1}^\dagger k_{t+1} \cdot \mathbf{x}_{t+1} (\mathbf{I} + k_{t+1} \cdot \mathbf{x}_{t+1}^T \eta_{t+1}^\dagger \mathbf{x}_{t+1})^{-1} \mathbf{x}_{t+1}^T \eta_{t+1}^\dagger$$

where $\eta_{t+1}^\dagger = (\eta_0^{-1} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^T)^{-1}$ is maintained as an intermediate variable for computation. This OL process of using $-kF$ -Bayes is still non-replay, can resist catastrophic forgetting, and adaptively determine the sensitive k factor.

Proof: Please see Appendix A.6

Remark 5. Given $\mathcal{B} = \{\mathbf{x}_i\}_{i=1}^b$ containing b samples, (66) is enhanced to:

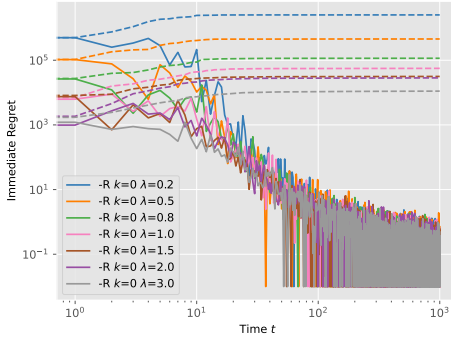
$$\begin{aligned} k_{t+1} &= \kappa \cdot \left(\frac{\text{trace}[(\mathcal{B}_{t+1} \eta_t \mathcal{B}_{t+1}^T + \sigma \mathbf{I})^{-1}]}{b} \right)^{-1} \\ k_t &= \kappa \cdot \left(\frac{\text{trace}[(\mathcal{B}_t \eta_t \mathcal{B}_t^T + \sigma \mathbf{I})^{-1}]}{b} \right)^{-1} \end{aligned} \quad (22)$$

4 EXPERIMENT

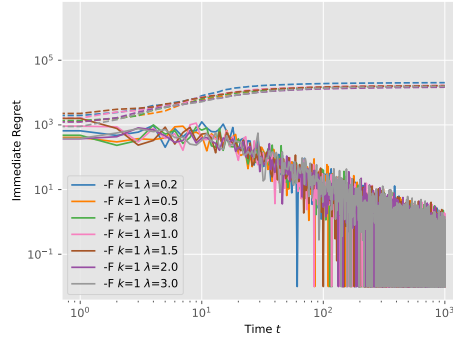
We first examined online learners with different regularization styles, namely -R, -F, $-kF$, and $-kF$ -Bayes in numerical simulations. Then we integrated these online learners into Randomized Neural Networks (Randomized NN) (Li & Zeng, 2023; Li et al., 2024) and conducted experiments in CIL/OTCIL scenarios, including tabular and image datasets. Based on the results, we show the efficacy of $-kF$ and $-kF$ -Bayes, and the great potential of Randomized learners equipped with them.

4.1 CASE 1: NUMERICAL SIMULATION

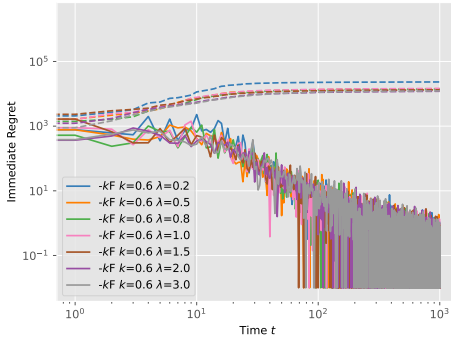
The numerical scenario was generated by the setups: $T=1000$, $d = 12$, $\varepsilon_t \sim \mathcal{N}(1, 2)$, $\mathbf{x}_t \sim \mathcal{N}_d(\mathbf{6I}, \Sigma)$, $\lambda \in \{0.2, 0.5, 0.8, 1.0, 1.5, 2.0, 3.0\}$, $k \in \{0.2, 0.4, 0.6, 0.8\}$, and the regularization styles included: -R, -F, $-kF$, $-kF$ -Bayes. We repeated the learning processes of the online learners 20 times using different random seeds for each λ . Results would be used to compute the mean and standard deviation of online regrets relative to Oracle learner (defined in (2)) in multiple trials. We assumed the learners had no prior knowledge before data came and started from $\theta = \mathbf{0}$.



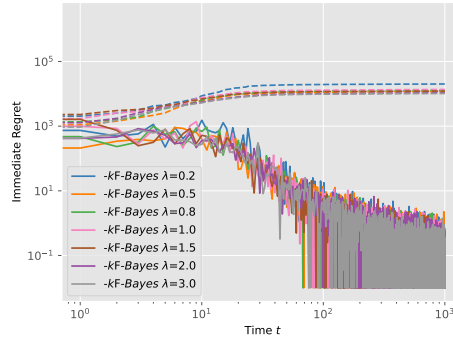
(a) Relative regrets of online learner using -R.



(b) Relative regrets of online learner using -F.



(c) Relative regrets of online learner using $-0.6F$.



(d) Relative regrets of online learner using $-kF$ -Bayes.

Figure 1: Relative immediate regrets (shown in solid lines) and cumulative regrets (shown in dashed lines) of online linear regression using different regularization styles. The shadows indicate ranges within one standard deviation. $-0.6F$ represents the best performance in all tested $-kF$, one can refer to Appendix A.7 for other results.

As shown in Figure 1, learners using -R, also regarded as $-0F$ style, suffer from low stability and high immediate regrets at all λ s. Although it converges quickly at a later stage, it causes high cumulative regrets in OL process. -F (i.e. $-1.0F$) style is more stable to -R and the cumulative error decreases significantly. Here we roughly selected several points for studying k 's effects. The best $-0.6F$ style slightly reduces cumulative regrets and responds faster than -F. The effect is not obvious because such a k exists but is difficult to determine and the same k is poorly adapted to different

environments, i.e., it does not perform as expected for all λ s. Appendix A.7 also shows that k is sensitive and an improper k results in a large loss. How to set a proper k is important because the entire dataset is unavailable and real-time searching is unreasonable in practice. The $-kF$ -Bayes achieves rapid decreases of immediate regrets and lower cumulative loss with the k s adaptively set in OL. The advantages of $-kF$ -Bayes will become more apparent in difficult tasks. The k variation curves are shown in Figure 2

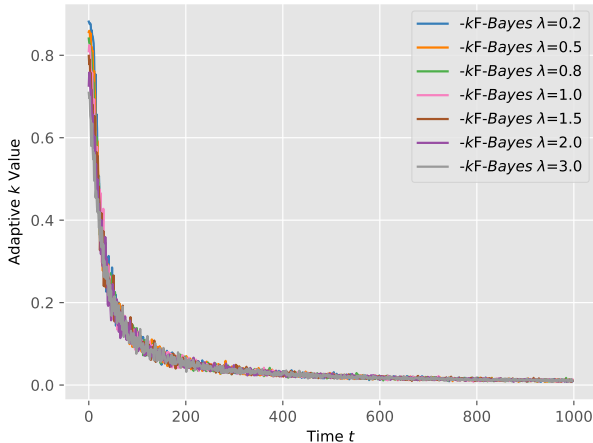


Figure 2: k variation curves of online learners with $-kF$ -Bayes in numerical simulations.

4.2 CASE 2: TABULAR DATASET

In order to examine linear regression learners using various regularization styles in more complex tasks, we used $-kF$ and $-kF$ -Bayes algorithms to reconstruct a state-of-the-art Randomized NN, the ensemble deep random vector functional link network (edRVFL) (Shi et al., 2021), and adapted edRVFL- \mathcal{A} to CIL/OTCIL scenarios, where \mathcal{A} denotes varied styles. The datasets are listed in Table 1 and each one was treated as a single class incremental problem, namely learning on $[\text{Dataset}] - |\mathcal{Y}|/|\mathcal{Y}|$ task stream.

Main methods involved in this experiment were as follows. (1) SMAC3 (Lindauer et al., 2022): we used it to configure well-behaved hyper-parameters (HP) within the same searching space for algorithms. (2) edRVFL (Shi et al., 2021): the original offline edRVFL network. (3) edRVFL-R, edRVFL- kF and edRVFL- kF -Bayes: our built methods for comparisons. We assumed they had no prior knowledge or update on any task, and always did one-pass training. We use N and L to denote edRVFL’s numbers of layers and nodes per layer respectively. (4) EWC (Kirkpatrick et al., 2017): the backbone of EWC for CIL was a BP-based MLP. (5) CRNet-I (Li & Zeng, 2023): a novel CIL network based on Randomized NN and EWC. (6) DYSON (He et al., 2024): an OTCIL method using compute-and-align paradigm. Some assess metrics containing ACC , BWT , and FWT are list in Appendix A.8 **Evaluation metrics**.

During SMAC3 operation, every dataset was randomly partitioned into 4-folds for consistent results, containing 60% for training, 15% for validation, and 25% for testing in each fold. We set trial times to 200. In one trial, the algorithm configured by SMAC3 would be tested on 4-folds separately using all different random seeds on task order and network initialization. The average $acc.(T)$ on 4-folds was set to be the cost of incumbents which guided SMAC3 to optimize. Note the k of $-kF$ was optimized while was adaptive in $-kF$ -Bayes.

Testset performance of the above algorithms with HP and structures optimized by SMAC3 is shown in Table 2 in Appendix A.8. Metrics are average $acc.(T)\%$ (for OTCIL), average $ACC(Q)\%$ (for CIL), and $std.\%$ of 4-fold trials. The results show that: (1) edRVFL- kF and edRVFL- kF -Bayes outperform other methods on most datasets; (2) although edRVFL- kF has favorable performance, it relies on fatal HP k value optimized by SMAC3. Searching methods (e.g. PSO, grid searching) are

usually time-consuming and difficult to optimize $-kF$ synchronously in CIL, while the $-kF$ -Bayes is ready to be deployed and achieves impressive results (even a little worse than $-kF$ occasionally); (3) the $-kF$ style using forward unsupervised information is not worse than $-R$ and when $k = 0$ they have similar results; (4) BP-based EWC and DYSON suffer from larger loss, and CRNet also shows undesirable capability, especially on longer task stream (e.g. letters, plant margin). (5) our methods are noticeably more stable as indicated by the low stds, and are better suited to this situation because no learning dissipation in theory. These conclusions also encourage the extension of our methods to the representative learning on image features transformed by pre-trained models (PTM).

We offer ACC curves on letters-26/26 as shown in Fig. 3. The optimized HP and structures were still employed. Our methods were set to no prior knowledge and start by $\theta_0 = 0$ and diagonal η_0 during OTCIL, which resulted in accuracy lags at $q \leq 5$. However, the hysteresis of edRVFL- kF -Bayes is slightest compared to using $-R$ or $-kF$, and all proposed methods arrive at the expected accuracy as more tasks are learned. EWC and CRNets yield increasing loss on long task stream, and DYSON gets around 20% more accuracy because of PTM and replay system. Note the optimization process like using SMAC3 is almost necessary for $-kF$ desired performance as its k requests careful setting, while is needless for $-kF$ -Bayes.

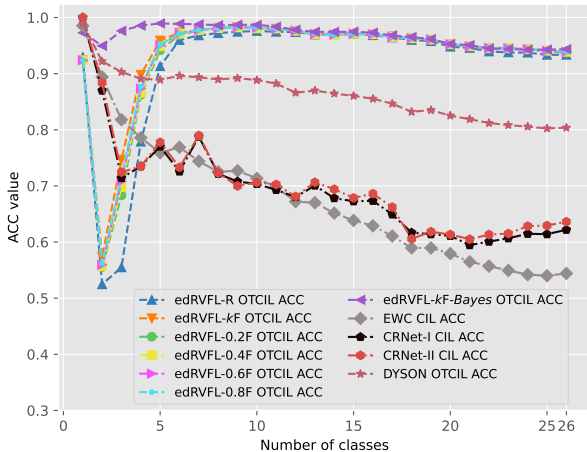


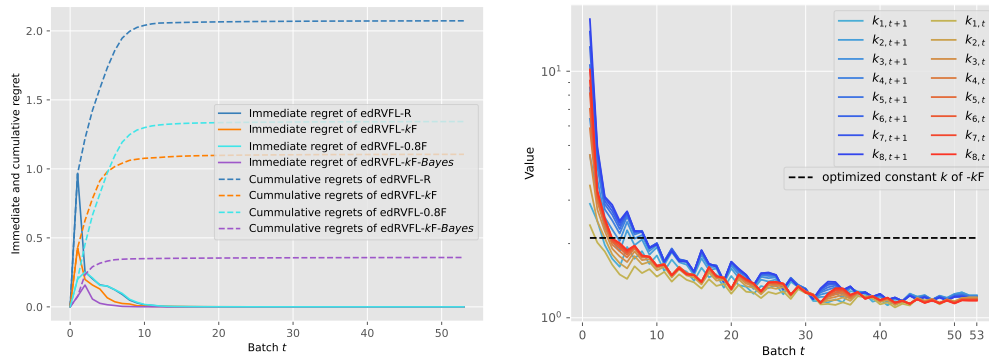
Figure 3: ACC curves of methods in Table 2 (in Appendix A.8) on letters-26/26 in CIL. The edRVFL-R, edRVFL- kF , and edRVFL- kF -Bayes learned one class in two batches without boundary.

The incurred loss during OTCIL is shown in Fig. 4 (a), which demonstrates the superiority of edRVFL- kF -Bayes in terms of accuracy, dynamic response, and regrets. The k s in $-kF$ were unreasonably identical in layers and time-consuming to tune, which was avoided and ready-to-use in $-kF$ -Bayes. The dynamic adaptive process of k in $-kF$ -Bayes is shown in Fig. 4 (b), and it pays more attention to a future task compared to past ones. The k of $-kF$ is sensitive and affected by multiple facts, and its being in the range of the maximum and minimum k_t of $-kF$ -Bayes is interesting, which highlights SMAC3’s choice and its effectiveness.

4.3 CASE 3: CIFAR IMAGE DATASET

For the CIFAR-100 dataset, we split it into CIFAR-100/10 to implement experiments. We studied selected methods, such as EWC (Kirkpatrick et al., 2017), CRNet (Li & Zeng, 2023) for the regularization-based group, RanPAC (McDonnell et al., 2024), NICE (Gurbuz et al., 2024), DYSON He et al. (2024) for the model-based methods, and GEM (Lopez-Paz & Ranzato, 2017), GSS (Aljundi et al., 2019) for the replay-based branch. A standard resnet-56 was employed as PTM and used to extract features from CIFAR-100, which inherited the setup in (Li & Zeng, 2023) for fair comparison, and to allow benchmarks to learn image representations. We offered the same PTM for all algorithms except when they are already equipped with one, such as for NICE and DYSON. After resnet-56, the PTM $\mathcal{F}(\cdot)$ was followed in our proposed methods for enhancement. The task

432
433
434
435
436
437
438
439
440
441
442
443
444



445
446
447
448
449
450

(a) The letter testset immediate and cumulative (b) $k_{l,t+1}$ and $k_{l,t}$ variation curves of $-kF$ -Bayes in OT-CIL. X-axis denotes batch, Y-axis is the logarithmic $Bayes$ in OT-CIL. $-kF$ -Bayes has the minimum loss. value. Trends illustrate the effect of maintaining the bal- X-axis denotes batch, and Y-axis is regret value. ance of penalties in changing optimization targets.

Figure 4: The regrets on testset in OT-CIL and adaptive k variation curves of edRVFL- kF -Bayes.

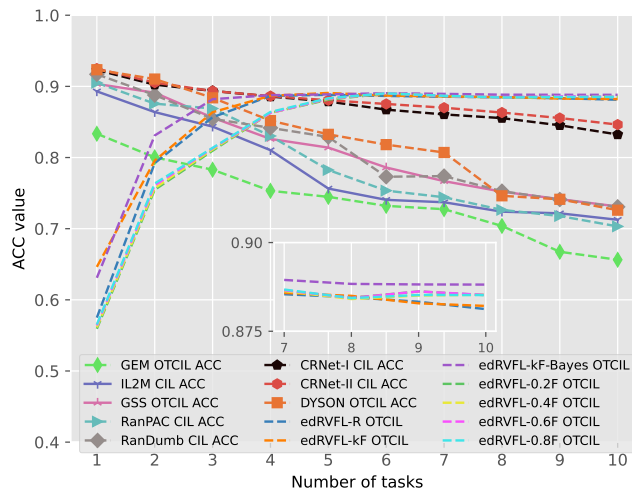
451
452
453
454
455

order was randomly sorted and baselines were tested 10 times. No task boundary was given to OT-CIL methods. Our methods still learned every task in one-pass, had no revisits, and we abode by severe no prior update before the data arrival.

456
457
458

The task order and choice of classes in CIFAR-100/10 were randomized. From Figure 5 and Figure 6, the proposed edRVFL- kF -Bayes still maintains an obvious advantage in performance.

460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476



477
478
479

Figure 5: ACC curves of methods on CIFAR-100/10 in CIL. Each task contained 10 classes. The edRVFL-R, edRVFL- kF , and edRVFL- kF -Bayes learned tasks in OT-CIL. X-axis is locally enlarged in the inlaid subfigures.

480
481
482
483

5 CONCLUSION

484
485

In this study, we introduce *Variable Forward regularization (-kF)* as an enhancement to the existing *Forward regularization (-F)* for online linear learners. Our findings indicate that -F, while theoretically promising, underperforms in practice due to inappropriate intervention penalties and

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

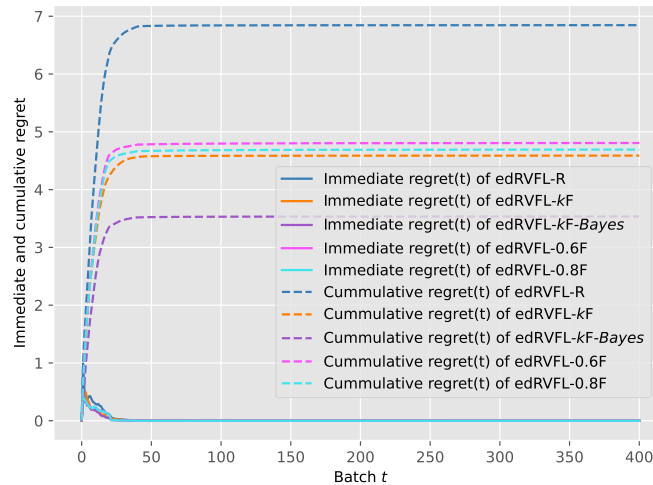


Figure 6: The CIFAR-100/10 testset immediate $regret(t)$ curves of edRVFL using -R, -kF, and -kF-Bayes strategies in OTCIL process. -kF-Bayes has the minimum loss in the whole process. X-axis denotes batch number, and Y-axis is regret and cumulative regret.

challenges posed by non-i.i.d nature in online learning. By modulating the intensity of -F with a variable k , -kF achieves improved stability and lower relative regret bounds, surpassing both -F and canonical *Ridge regularization* (-R). Additionally, we developed -kF-Bayes to dynamically adapt the penalty based on Bayesian principles, further addressing the instability issues. Experimental results in class incremental learning scenarios confirm the efficacy of our methods, highlighting their potential for reducing forgetting and enhancing performance in online learning tasks. Our future studies will focus on concrete regret bounds of -kF under stochastic setups.

REFERENCES

- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *Advances in neural information processing systems*, 32, 2019.
- Katy S Azoury and Manfred K Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. *Machine learning*, 43:211–246, 2001.
- Thomas Kleine Buening, Aadirupa Saha, Christos Dimitrakakis, and Haifeng Xu. Bandits meet mechanism design to combat clickbait in online recommendation. In *The Twelfth International Conference on Learning Representations*.
- Somnath Basu Roy Chowdhury, Nicholas Monath, Ahmad Beirami, Rahul Kidambi, Avinava Dubey, Amr Ahmed, and Snigdha Chaturvedi. Enhancing group fairness in online settings using oblique decision forests. *arXiv preprint arXiv:2310.11401*, 2023.
- Riccardo Della Vecchia and Debraj Basu. Online instrumental variable regression: Regret analysis and bandit feedback. *arXiv preprint arXiv:2302.09357*, 2023.
- Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999.
- Mustafa Burak Gurbuz, Jean Michael Moorman, and Constantine Dovrolis. Nice: Neurogenesis inspired contextual encoding for replay-free class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23659–23669, 2024.
- Yuhang He, Yingjie Chen, Yuhan Jin, Songlin Dong, Xing Wei, and Yihong Gong. Dyson: Dynamic feature space self-organization for online task-free class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23741–23751, 2024.

- 540 Lu Jiang, Yanan Xiao, Xinxin Zhao, Yuanbo Xu, Shuli Hu, Pengyang Wang, and Minghao Yin.
541 Hierarchical reinforcement learning on multi-channel hypergraph neural network for course rec-
542 ommendation. In *Proceedings of the Thirty-Third International Joint Conference on Artificial*
543 *Intelligence (IJCAI-24)*, 2024.
- 544 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
545 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcom-
546 ing catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*,
547 114(13):3521–3526, 2017.
- 548 Depeng Li and Zhigang Zeng. Crnet: A fast continual learning framework with random theory.
549 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):10731–10744, 2023.
- 550 Depeng Li, Tianqi Wang, Junwei Chen, Wei Dai, and Zhigang Zeng. Harnessing neural unit dynam-
551 ics for effective and scalable class-incremental learning. In *Forty-first International Conference*
552 *on Machine Learning*, 2024.
- 553 Marius Lindauer, Katharina Eggenberger, Matthias Feurer, André Biedenkapp, Difan Deng, Car-
554 olin Benjamins, Tim Ruhkopf, René Sass, and Frank Hutter. Smac3: A versatile bayesian opti-
555 mization package for hyperparameter optimization. *The Journal of Machine Learning Research*,
556 23(1):2475–2483, 2022.
- 557 David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning.
558 *Advances in neural information processing systems*, 30, 2017.
- 559 Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasnejad, and Anton van den Hengel.
560 Ranpac: Random projections and pre-trained models for continual learning. *Advances in Neural*
561 *Information Processing Systems*, 36, 2024.
- 562 Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye
563 Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for contin-
564 ual learning. In *International conference on machine learning*, pp. 4528–4537. PMLR, 2018.
- 565 Qiushi Shi, Rakesh Katuwal, Ponnuthurai N Suganthan, and Muhammad Tanveer. Random vector
566 functional link neural network based ensemble deep learning. *Pattern Recognition*, 117:107978,
567 2021.
- 568 Sam Toyer, Olivia Watkins, Ethan Adrian Mendes, Justin Svegliato, Luke Bailey, Tiffany Wang,
569 Isaac Ong, Karim Elmaaroufi, Pieter Abbeel, Trevor Darrell, et al. Tensor trust: Interpretable
570 prompt injection attacks from an online game. In *The Twelfth International Conference on Learn-*
571 *ing Representations*, 2024.
- 572 Volodya Vovk. Competitive on-line statistics. *International Statistical Review*, 69(2):213–248, 2001.
- 573 Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence.
574 In *International conference on machine learning*, pp. 3987–3995. PMLR, 2017.
- 575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

A APPENDIX

A.1 THEOREM 2 PROOF

Proof: (7) can be expanded to:

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} U_t(\boldsymbol{\theta}) - U_t(\boldsymbol{\theta}_t) + \mathcal{L}_t(\boldsymbol{\theta}) - (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T \nabla U_t(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}) \quad (23)$$

The $\nabla U_t(\boldsymbol{\theta}_t) = 0$ because of the latest convex optimization. (23) can be simplified to:

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} U_{t+1}(\boldsymbol{\theta}) - U_t(\boldsymbol{\theta}_t) = \operatorname{argmin}_{\boldsymbol{\theta}} \Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathcal{L}_{1..t}(\boldsymbol{\theta}) + k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - c. \quad (24)$$

Solutions to both models (i.e. online learner (7) and offline expert (24)) remain the same at each time point t if designating the last target as Bregman divergence function. This also suggests that $-kF$ suffers no learning dissipation compared to the offline expert, because they have similar optimum at each step.

Proof finished. \square

A.2 THEOREM 3 PROOF

Proof: Concretizing Theorem 1 with the **Setup**. For $0 \leq t \leq T$ we have:

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} U_{t+1}(\boldsymbol{\theta}) \quad (25)$$

$$U_{t+1}(\boldsymbol{\theta}) = \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \boldsymbol{\eta}_0^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \sum_{i=1}^t \frac{1}{2} \|\mathbf{x}_i^T \boldsymbol{\theta} - y_i\|_2^2 + \frac{k}{2} \|\mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2$$

Convert (25) into the form of (7) described by Theorem 2 to avoid the retrospective retraining, and derive variable learning rate via Lemma 2:

$$U_t(\boldsymbol{\theta}) + \frac{1}{2} \|\mathbf{x}_t^T \boldsymbol{\theta} - y_t\|_2^2 + \frac{k}{2} \|\mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2 - \frac{k}{2} \|\mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2 \quad (26)$$

$$= U_0(\boldsymbol{\theta}) - U_0(\boldsymbol{\theta}_0) - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \nabla U_0(\boldsymbol{\theta}_0) + \sum_{i=1}^t \frac{1}{2} \|\mathbf{x}_i^T \boldsymbol{\theta} - y_i\|_2^2 + \frac{k}{2} \|\mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2$$

$$\Rightarrow U_t(\boldsymbol{\theta}) - U_0(\boldsymbol{\theta}) = \sum_{i=1}^{t-1} \frac{1}{2} \|\mathbf{x}_i^T \boldsymbol{\theta} - y_i\|_2^2 + \frac{k}{2} \|\mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2 - \frac{1}{2} \boldsymbol{\theta}_0^T \boldsymbol{\eta}_0^{-1} \boldsymbol{\theta}_0 - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \boldsymbol{\eta}_0^{-1} \boldsymbol{\theta}_0$$

$$\Rightarrow U_t(\boldsymbol{\theta}) - U_0(\boldsymbol{\theta}) = \sum_{i=1}^{t-1} \frac{1}{2} \|\mathbf{x}_i^T \boldsymbol{\theta} - y_i\|_2^2 + \frac{k}{2} \|\mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2 + \frac{1}{2} \boldsymbol{\theta}_0^T \boldsymbol{\eta}_0^{-1} \boldsymbol{\theta}_0 - \boldsymbol{\theta}^T \boldsymbol{\eta}_0^{-1} \boldsymbol{\theta}_0$$

$$\Rightarrow U_t(\boldsymbol{\theta}) - U_0(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^T \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T \boldsymbol{\theta} + \frac{1}{2} \sum_{i=1}^{t-1} y_i^2 - \sum_{i=1}^{t-1} y_i \mathbf{x}_i^T \boldsymbol{\theta}$$

$$- \boldsymbol{\theta}^T \boldsymbol{\eta}_0^{-1} \boldsymbol{\theta}_0 + \frac{k}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{x}_t \mathbf{x}_t^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \frac{1}{2} \boldsymbol{\theta}_0^T \boldsymbol{\eta}_0^{-1} \boldsymbol{\theta}_0$$

Based on Lemma 2, (26) can be simplified to:

$$U_t(\boldsymbol{\theta}) - (U_0(\boldsymbol{\theta}) + \frac{1}{2} \boldsymbol{\theta}^T \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T \boldsymbol{\theta} + \frac{k}{2} \boldsymbol{\theta}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}) = \boldsymbol{\omega}^T \boldsymbol{\theta} + c. \quad (27)$$

where $\boldsymbol{\omega} \in \mathbb{R}^d$. If set $V_t = \frac{1}{2} \boldsymbol{\theta}^T (\sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T + k \cdot \mathbf{x}_t \mathbf{x}_t^T) \boldsymbol{\theta}$, based on (27) and Lemma 1 we have:

$$\Delta_{U_t}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \Delta_{U_0 + V_t}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) + \Delta_{V_t}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) \quad (28)$$

$$= \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T [\boldsymbol{\eta}_0^{-1} + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T + k \cdot \mathbf{x}_t \mathbf{x}_t^T] (\boldsymbol{\theta} - \boldsymbol{\theta}_t)$$

(28) shows that it alleviates the previous replay with all that knowledge absorbed in Bregman divergence. The stepwise solution to (28) in OL will change with $\boldsymbol{\eta}_t^{-1} = \boldsymbol{\eta}_0^{-1} + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T + k \cdot \mathbf{x}_t \mathbf{x}_t^T$, also termed as the variable learning rate. The essence of OL algorithm using $-kF$ is also to give time-varying optimization targets to guide network optimum (a posterior from the Bayesian view) updates on task streams. This optimization objective in Theorem 2 can be expressed further as:

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} \Delta_{U_t}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) + \mathcal{L}_t(\boldsymbol{\theta}) + k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}) \quad (29)$$

$$= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T [\boldsymbol{\eta}_0^{-1} + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T + k \cdot \mathbf{x}_t \mathbf{x}_t^T] (\boldsymbol{\theta} - \boldsymbol{\theta}_t)$$

$$+ \frac{1}{2} \|\mathbf{x}_t^T \boldsymbol{\theta} - y_t\|_2^2 + \frac{k}{2} \|\mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2 - \frac{k}{2} \|\mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2$$

(29) can be transformed to the following constrained optimization problem:

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \boldsymbol{\theta}_{t+1}} \quad & \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T \boldsymbol{\eta}_t^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2}\|\xi_1\|_2^2 + \frac{k}{2}\|\xi_2\|_2^2 - \frac{k}{2}\|\xi_3\|_2^2 \\ \text{s.t.} \quad & \mathbf{x}_t^T \boldsymbol{\theta} - y_t = \xi_1; \mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = \xi_2; \mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = \xi_3, \forall t \end{aligned} \quad (30)$$

The Lagrangian function of problem (30) is:

$$\begin{aligned} \ell(\boldsymbol{\theta}, \xi_{\{1,2,3\}}, \mu_{\{1,2,3\}}) = & \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T \boldsymbol{\eta}_t^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2}\|\xi_1\|_2^2 + \frac{k}{2}\|\xi_2\|_2^2 - \frac{k}{2}\|\xi_3\|_2^2 \\ & + \mu_1(\mathbf{x}_t^T \boldsymbol{\theta} - y_t - \xi_1) + \mu_2(\mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) - \xi_2) + \mu_3(\mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) - \xi_3) \end{aligned} \quad (31)$$

The Lagrangian function (31) can be tackled through Karush-Kuhn-Tucker (KKT) conditions, which can be built into the following formulation:

$$\begin{aligned} \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \boldsymbol{\theta}} = 0 & \Rightarrow \boldsymbol{\eta}_t^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \mu_1 \mathbf{x}_t + \mu_2 \mathbf{x}_{t+1} + \mu_3 \mathbf{x}_t = 0 \\ \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \xi_1} = 0 & \Rightarrow \xi_1 = \mu_1 \\ \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \xi_2} = 0 & \Rightarrow k \xi_2 = \mu_2 \\ \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \xi_3} = 0 & \Rightarrow -k \xi_3 = \mu_3 \\ \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \mu_1} = 0 & \Rightarrow \mathbf{x}_t^T \boldsymbol{\theta} - y_t = \xi_1 \\ \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \mu_2} = 0 & \Rightarrow \mathbf{x}_{t+1}^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = \xi_2 \\ \frac{\partial \ell(\boldsymbol{\theta}, \xi, \mu)}{\partial \mu_3} = 0 & \Rightarrow \mathbf{x}_t^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = \xi_3 \end{aligned} \quad (32)$$

Based on (32), we can obtain the recursive updating policy between $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}_{t+1}$ as follows:

$$\begin{aligned} \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \boldsymbol{\eta}_{t+1}[(\mathbf{x}_t \mathbf{x}_t^T + k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T - k \cdot \mathbf{x}_t \mathbf{x}_t^T) \boldsymbol{\theta}_t - \mathbf{x}_t y_t] \\ + k \cdot \boldsymbol{\eta}_{t+1}(\mathbf{x}_{t+1} \mathbf{x}_{t+1}^T - \mathbf{x}_t \mathbf{x}_t^T) \boldsymbol{\theta}_0 \end{aligned} \quad (33)$$

Without loss of generality, given $\boldsymbol{\theta}_0 = 0$, to allow the learner to start from a blank model (no prior knowledge), (33) is rewritten as:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \boldsymbol{\eta}_{t+1}[(\mathbf{x}_t \mathbf{x}_t^T + k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T - k \cdot \mathbf{x}_t \mathbf{x}_t^T) \boldsymbol{\theta}_t - \mathbf{x}_t y_t] \quad (34)$$

The variable learning rate $\boldsymbol{\eta}^{-1}$ is also updated step-wise by using Sherman–Morrison–Woodbury law, as shown in (9).

Proof finished. \square

A.3 REMARK 2 PROOF

Proof: (34) can be rewritten as:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \boldsymbol{\eta}_{t+1}[(\boldsymbol{\eta}_{t+1}^{-1} - \boldsymbol{\eta}_t^{-1}) \boldsymbol{\theta}_t - \mathbf{x}_t y_t] \Rightarrow (12) \quad (35)$$

(12) can be rewritten as:

$$\begin{aligned} \boldsymbol{\eta}_t \boldsymbol{\eta}_{t+1}^{-1} \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \boldsymbol{\eta}_t \mathbf{x}_t y_t \\ \boldsymbol{\eta}_t (\boldsymbol{\eta}_t^{-1} + (1-k) \cdot \mathbf{x}_t \mathbf{x}_t^T + k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T) \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \boldsymbol{\eta}_t \mathbf{x}_t y_t \Rightarrow (13) \end{aligned} \quad (36)$$

Take the derivation of 25 to calculate $\boldsymbol{\theta}_{T+1}$ when $t = T$:

$$\left. \frac{\partial U_{T+1}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{T+1}} = 0 \Rightarrow \boldsymbol{\theta}_{T+1} = \boldsymbol{\eta}_{T+1} \sum_{i=1}^T \mathbf{x}_i y_i \Rightarrow (14) \quad (37)$$

Proof finished. \square

A.4 THEOREM 4 PROOF

Proof: For $0 \leq t \leq T$, we expand the divergence $\Delta_{U_{t+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{t+1})$:

$$\begin{aligned}\Delta_{U_{t+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{t+1}) &= U_{t+1}(\boldsymbol{\theta}) - U_{t+1}(\boldsymbol{\theta}_{t+1}) - (\boldsymbol{\theta} - \boldsymbol{\theta}_{t+1})^T \nabla_{\boldsymbol{\theta}} U_{t+1}(\boldsymbol{\theta}_{t+1}) \\ &= U_{t+1}(\boldsymbol{\theta}) - U_{t+1}(\boldsymbol{\theta}_{t+1})\end{aligned}\quad (38)$$

According to Theorem 2, we have:

$$U_{t+1}(\boldsymbol{\theta}) = U_t(\boldsymbol{\theta}) + \mathcal{L}_t(\boldsymbol{\theta}) + k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta})\quad (39)$$

Substitute (38) into (39), for $1 \leq t \leq T$, we get:

$$\mathcal{L}_t(\boldsymbol{\theta}) = \Delta_{U_{t+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{t+1}) + U_{t+1}(\boldsymbol{\theta}_{t+1}) - U_t(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta})\quad (40)$$

Let $\boldsymbol{\theta} = \boldsymbol{\theta}_t$ in (40), we obtain:

$$\mathcal{L}_t(\boldsymbol{\theta}_t) = \Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) + U_{t+1}(\boldsymbol{\theta}_{t+1}) - U_t(\boldsymbol{\theta}_t) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)\quad (41)$$

Subtract (40) from (41), we have:

$$\begin{aligned}\mathcal{L}_t(\boldsymbol{\theta}_t) - \mathcal{L}_t(\boldsymbol{\theta}) &= \Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t) \\ &\quad + k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}) - \Delta_{U_{t+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{t+1}) - U_t(\boldsymbol{\theta}_t) + U_t(\boldsymbol{\theta})\end{aligned}\quad (42)$$

Substitute (38) into (42) again, we obtain:

$$\begin{aligned}\mathcal{L}_t(\boldsymbol{\theta}_t) - \mathcal{L}_t(\boldsymbol{\theta}) &= \Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t) \\ &\quad + k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}) - \Delta_{U_{t+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{t+1}) + \Delta_{U_t}(\boldsymbol{\theta}, \boldsymbol{\theta}_t)\end{aligned}\quad (43)$$

Integrate both sides of (43) over $1 \leq t \leq T$ and subtract $\Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0)$, we have:

$$\begin{aligned}\sum_{t=1}^T \mathcal{L}_t(\boldsymbol{\theta}_t) - (\Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \sum_{t=1}^T \mathcal{L}_t(\boldsymbol{\theta})) &= \sum_{t=1}^T [\Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)] \\ &\quad + k \cdot \hat{\mathcal{L}}_{T+1}(\boldsymbol{\theta}) - k \cdot \hat{\mathcal{L}}_1(\boldsymbol{\theta}) - \Delta_{U_{T+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{T+1}) + \Delta_{U_1}(\boldsymbol{\theta}, \boldsymbol{\theta}_1) - \Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0)\end{aligned}\quad (44)$$

Note that the last few terms in (44) are iteratively canceled out.

Let $t = 0$ in (39) and $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_0 = 0$, we have $U_1(\boldsymbol{\theta}) = U_0(\boldsymbol{\theta}) + k \cdot \hat{\mathcal{L}}_1(\boldsymbol{\theta})$.

$$\Delta_{U_1}(\boldsymbol{\theta}, \boldsymbol{\theta}_1) - \Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) = U_1(\boldsymbol{\theta}) - U_1(\boldsymbol{\theta}_1) - U_0(\boldsymbol{\theta}) + U_0(\boldsymbol{\theta}_0) = k \cdot \hat{\mathcal{L}}_1(\boldsymbol{\theta})\quad (45)$$

Finally, we can obtain:

$$\begin{aligned}\sum_{t=1}^T \mathcal{L}_t(\boldsymbol{\theta}_t) - (\Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \sum_{t=1}^T \mathcal{L}_t(\boldsymbol{\theta})) \\ = \sum_{t=1}^T [\Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)] + k \cdot \hat{\mathcal{L}}_{T+1}(\boldsymbol{\theta}) - \Delta_{U_{T+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{T+1})\end{aligned}\quad (46)$$

Proof finished. \square

A.5 THEOREM 5 PROOF

Proof: Use Lemma 1 and Setups to concrete the right part of (15) in Theorem 4:

$$\begin{aligned}\sum_{t=1}^T [\Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)] + k \cdot \hat{\mathcal{L}}_{T+1}(\boldsymbol{\theta}) - \Delta_{U_{T+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{T+1}) \\ = \sum_{t=1}^T \left[\frac{1}{2} (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1})^T \boldsymbol{\eta}_{t+1}^{-1} (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}) - \frac{k}{2} \|\mathbf{x}_{t+1}^T \boldsymbol{\theta}_t\|_2^2 + \frac{k}{2} \|\mathbf{x}_t^T \boldsymbol{\theta}_t\|_2^2 \right] \\ - \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{T+1})^T \boldsymbol{\eta}_{T+1}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_{T+1}) + \frac{k}{2} \|\mathbf{x}_{T+1}^T \boldsymbol{\theta}\|_2^2\end{aligned}\quad (47)$$

756 Substitute (34) (12) (13) into the first term of (47):
 757
 758
 759
 760
 761
 762

$$\begin{aligned}
 & \frac{1}{2}(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1})^T \boldsymbol{\eta}_{t+1}^{-1} (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}) & (48) \\
 &= \frac{1}{2}(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1})^T (k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_t + (1-k) \cdot \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - \mathbf{x}_t y_t) \\
 &= \frac{1}{2}(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1})^T (-k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} + (1-k) \cdot \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - \mathbf{x}_t y_t) \\
 &+ \frac{1}{2}(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1})^T (k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T) (\boldsymbol{\theta}_t + \boldsymbol{\theta}_{t+1}) \\
 &= \frac{1}{2} [(1-k) \cdot \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_{t+1} + k \cdot \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} - \boldsymbol{\eta}_t \mathbf{x}_t y_t]^T \\
 &\times [(1-k) \cdot \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} - \mathbf{x}_t y_t] + \frac{1}{2} (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1})^T (k \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T) (\boldsymbol{\theta}_t + \boldsymbol{\theta}_{t+1}) \\
 &= \frac{1}{2} [(1-k)^2 \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} \\
 &- (1-k) \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t y_t + k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t \\
 &- k^2 \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} - k \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_t y_t \\
 &- (1-k) \cdot y_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t + k \cdot y_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} + y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t] \\
 &+ k \cdot \frac{1}{2} \boldsymbol{\theta}_t^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_t - k \cdot \frac{1}{2} \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}
 \end{aligned}$$

763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790 Substitute (14) into the last two terms of (47) and let $\boldsymbol{\eta}_* = (\boldsymbol{\eta}_0^{-1} + \sum_{i=1}^T \mathbf{x}_i \mathbf{x}_i^T)^{-1} \Rightarrow \boldsymbol{\theta}^{off} -$
 791 $\boldsymbol{\theta}_{T+1} = k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1}$, we obtain:
 792
 793
 794
 795
 796
 797

$$\begin{aligned}
 & -\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{T+1})^T \boldsymbol{\eta}_{T+1}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_{T+1}) + \frac{k}{2} \|\mathbf{x}_{T+1}^T \boldsymbol{\theta}\|_2^2 & (49) \\
 &= -\frac{1}{2} k^2 \cdot \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\eta}_* \boldsymbol{\eta}_{T+1}^{-1} \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1} \\
 &+ \frac{k}{2} \boldsymbol{\theta}_{T+1}^T (\mathbf{I} + k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T)^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T (\mathbf{I} + k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T) \boldsymbol{\theta}_{T+1} \\
 &= -\frac{1}{2} k^2 \cdot \boldsymbol{\theta}_{T+1}^T (\mathbf{x}_{T+1} \mathbf{x}_{T+1}^T)^2 \boldsymbol{\eta}_* (\mathbf{I} + k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T) \boldsymbol{\theta}_{T+1} \\
 &+ \frac{k}{2} \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T (\mathbf{I} + k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T)^2 \boldsymbol{\theta}_{T+1} \\
 &= \frac{k}{2} \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T (\mathbf{I} + k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T) \boldsymbol{\theta}_{T+1}
 \end{aligned}$$

Substitute (48) and (49) back into (47), we get:

$$\begin{aligned}
& \sum_{t=1}^T [\Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)] + k \cdot \hat{\mathcal{L}}_{T+1}(\boldsymbol{\theta}) - \Delta_{U_{T+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{T+1}) \quad (50) \\
&= \frac{1}{2} \sum_{t=1}^T [(1-k)^2 \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} \\
&\quad - (1-k) \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t y_t + k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t \\
&\quad - k^2 \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} - (1-k) \cdot y_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t + y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \\
&\quad + k \cdot \boldsymbol{\theta}_t^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - k \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}] + \frac{k}{2} \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T (\mathbf{I} + k \cdot \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T) \boldsymbol{\theta}_{T+1} \\
&= \frac{1}{2} \sum_{t=1}^T [(1-k)^2 \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} \\
&\quad - (1-k) \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t y_t + k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t \\
&\quad - k^2 \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} - (1-k) \cdot y_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t + y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t] \\
&\quad + \frac{k^2}{2} \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1}
\end{aligned}$$

Try to simplify (50) further. When $t = T$ and using Sherman–Morrison–Woodbury law, we have:

$$\begin{aligned}
& \frac{k^2}{2} \cdot \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\eta}_* \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1} - \frac{k^2}{2} \cdot \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\eta}_T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1} \quad (51) \\
&= \frac{k^2}{2} \cdot \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T [\boldsymbol{\eta}_* - \boldsymbol{\eta}_T] \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1} \\
&= \frac{k^2}{2[1 + (1-k) \mathbf{x}_T^T \boldsymbol{\eta}_T \mathbf{x}_T]} \cdot \boldsymbol{\theta}_{T+1}^T \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T [-(1-k) \boldsymbol{\eta}_T \mathbf{x}_T \mathbf{x}_T^T \boldsymbol{\eta}_T] \mathbf{x}_{T+1} \mathbf{x}_{T+1}^T \boldsymbol{\theta}_{T+1}
\end{aligned}$$

According to the property of the symmetric definite matrix, given $0 < k < 1$ and $\mathbf{x} \neq \mathbf{0}$, (51) < 0. So the initial two terms in (51) can be eliminated in (50) as we target the upper restrictions. The (50) can be rewritten as:

$$\begin{aligned}
& \sum_{t=1}^T [\Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)] + k \cdot \hat{\mathcal{L}}_{T+1}(\boldsymbol{\theta}) - \Delta_{U_{T+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{T+1}) \quad (52) \\
&= \frac{1}{2} \sum_{t=1}^T [(1-k)^2 \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} \\
&\quad - (1-k) \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t y_t + k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t \\
&\quad - (1-k) \cdot y_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t + y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t] - \frac{k^2}{2} \cdot \sum_{t=1}^{T-1} \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}
\end{aligned}$$

Given $\mathcal{X} \sim \mathcal{P}_{\mathcal{X}}$ and $\{\mathbf{x}_t\}_{t=1}^T$ respects the i.i.d premise,

$$\begin{aligned}
& \mathbb{E}[(1-k)^2 \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t + k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t] \quad (53) \\
&\leq \mathbb{E}[(1-k)^2 \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t + k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}] \\
&= \mathbb{E}[(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t]
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}[(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - (1-k) \cdot y_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t] \quad (54) \\
&= \mathbb{E}[(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t - (1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\theta}_t] \\
&= 0
\end{aligned}$$

Substitute (53) and (54) back into (52), we obtain:

$$\begin{aligned}
& \mathbb{E}[\sum_{t=1}^T [\Delta_{U_{t+1}}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}) - k \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}_t) + k \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}_t)] + k \cdot \hat{\mathcal{L}}_{T+1}(\boldsymbol{\theta}) - \Delta_{U_{T+1}}(\boldsymbol{\theta}, \boldsymbol{\theta}_{T+1})] \quad (55) \\
&= \mathbb{E}[\frac{1}{2} \sum_{t=1}^T [-k(1-k) \cdot \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1} - (1-k) \boldsymbol{\theta}_{t+1}^T \mathbf{x}_t \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{y}_t + y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t] \\
&\quad - \frac{k^2}{2} \cdot \sum_{t=1}^{T-1} \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}] \\
&= \mathbb{E}[\frac{1}{2} \sum_{t=1}^T [-(1-k) \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \mathbf{y}_t^2 + y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t] - \frac{k^2}{2} \cdot \sum_{t=1}^{T-1} \boldsymbol{\theta}_{t+1}^T \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}] \\
&\leq \mathbb{E}[\frac{k}{2} \sum_{t=1}^T y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t]
\end{aligned}$$

To prevent the model from potential attacks, in the adversarial setting, we assume $\mathbf{x}_t^T \boldsymbol{\theta}_t$ lies in $[-Y_m, Y_m]$ for -R style as shown in Lemma 4. Obviously, our general $-kF$ style removes this and avoids clipping or Y_m updates. To constrain (55) further, according to the Matrix Spectral Theorem, it can be proved that $e_{max}(\mathbf{A}) = \sup(\mathbf{p}^T \mathbf{A} \mathbf{p} \mid \|\mathbf{p}\|_2 = 1)$ is a convex function of \mathbf{A} , where e_{max} serves as the maximum eigenvalue of real-valued symmetric matrix \mathbf{A} , (55) can be transformed to:

$$\begin{aligned}
& \mathbb{E}[\frac{k}{2} \sum_{t=1}^T y_t^2 \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t] \quad (56) \\
&= \mathbb{E}[\frac{k}{2} \sum_{t=1}^T y_t^2 \|\mathbf{x}_t\|_2^2 \frac{\mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t}{\|\mathbf{x}_t\|_2 \|\boldsymbol{\eta}_t \mathbf{x}_t\|_2}] \\
&\leq \mathbb{E}[\frac{k}{2} Y_m^2 X_m^2 d \sum_{t=1}^T e_{max}(\boldsymbol{\eta}_t)] \\
&= \mathbb{E}[\frac{k}{2} Y_m^2 X_m^2 d \sum_{t=1}^T \inf[\frac{1}{e_{min}(\boldsymbol{\eta}_t^{-1})}]]
\end{aligned}$$

where $X_m = \max_{1 \leq t \leq T} \{\|\mathbf{x}_t\|_\infty\}$, $Y_m = \max_{1 \leq t \leq T} \{y_t\}$, and $\sup(e_{min}(\boldsymbol{\eta}_t^{-1}) - \boldsymbol{\eta}_0^{-1}) \leq X_m^2(t-1+k)$ because $\sum_i e_i(\boldsymbol{\eta}_t^{-1} - \boldsymbol{\eta}_0^{-1}) = \text{trace}(\boldsymbol{\eta}_t^{-1} - \boldsymbol{\eta}_0^{-1})$. Continue to simplify (56), we obtain:

$$\begin{aligned}
& \mathbb{E}[\frac{k}{2} Y_m^2 X_m^2 d \sum_{t=1}^T \inf[\frac{1}{e_{min}(\boldsymbol{\eta}_t^{-1})}]] \quad (57) \\
&\leq \mathbb{E}[\frac{k}{2} Y_m^2 X_m^2 d \int_0^T \frac{1}{\lambda + X_m^2(t-1+k)} dt] \\
&= \mathbb{E}[\frac{k}{2} Y_m^2 d \int_0^T \frac{1}{t + \frac{\lambda}{X_m^2} + k - 1}] \\
&= \frac{k}{2} Y_m^2 d \ln(1 + \frac{TX_m^2}{\lambda + (k-1) \cdot X_m^2})
\end{aligned}$$

Proof finished. \square

A.6 THEOREM 6 PROOF

Proof: To enable $-kF$ to adaptively mediate the Forward regularization intensity for robust performance in OL/CL, the hard constant k is replaced with soft variables. We will study the dynamic

updating process and introduce an algorithm for generating time-varying k values. (6) in Theorem 1 can be written as:

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &= \operatorname{argmin}_{\boldsymbol{\theta}} U_{t+1}(\boldsymbol{\theta}) \\ U_{t+1}(\boldsymbol{\theta}) &= \Delta_{U_0}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathcal{L}_{1..t}(\boldsymbol{\theta}) + k_{t+1} \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta})\end{aligned}\quad (58)$$

Based on (58), we represent Theorem 2 as follows:

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} \Delta_{U_t}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) + \mathcal{L}_t(\boldsymbol{\theta}) + k_{t+1} \cdot \hat{\mathcal{L}}_{t+1}(\boldsymbol{\theta}) - k_t \cdot \hat{\mathcal{L}}_t(\boldsymbol{\theta}) \quad (59)$$

Similar to the proof of Theorem 3, we can obtain:

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T [\boldsymbol{\eta}_0^{-1} + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T + k_t \cdot \mathbf{x}_t \mathbf{x}_t^T] (\boldsymbol{\theta} - \boldsymbol{\theta}_t) \\ &\quad + \frac{1}{2} \|\mathbf{x}_t^T \boldsymbol{\theta} - y_t\|_2^2 + \frac{k_{t+1}}{2} \cdot \|\mathbf{x}_{t+1}^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2 - \frac{k_t}{2} \cdot \|\mathbf{x}_t^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_2^2\end{aligned}\quad (60)$$

where the novel variable learning rate $\boldsymbol{\eta}_t^{-1} = \boldsymbol{\eta}_0^{-1} + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^T + k_t \cdot \mathbf{x}_t \mathbf{x}_t^T$. (34) is refreshed to:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \boldsymbol{\eta}_{t+1} [(\mathbf{x}_t \mathbf{x}_t^T + k_{t+1} \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T - k_t \cdot \mathbf{x}_t \mathbf{x}_t^T) \boldsymbol{\theta}_t - \mathbf{x}_t y_t] \quad (61)$$

where the novel variable learning rate $\boldsymbol{\eta}_{t+1}^{-1} = \boldsymbol{\eta}_0^{-1} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^T + k_{t+1} \cdot \mathbf{x}_{t+1} \mathbf{x}_{t+1}^T$.

From the perspective of Bayesian learning:

$$p(\boldsymbol{\theta}_{t+1} | y_1, \dots, y_t) \propto p(y_t | \boldsymbol{\theta}_{t+1}, y_1, \dots, y_{t-1}) \cdot p(\boldsymbol{\theta}_{t+1} | y_1, \dots, y_{t-1}), \quad (62)$$

the OL/CL can be regarded as a process of consecutively updating the posterior distribution of learnable weights based on current empirical risk and the prior distribution.

In (60), the prior Gaussian distribution of $\boldsymbol{\theta}$ can be:

$$p(\boldsymbol{\theta}_{t+1}) \sim \mathcal{N}(\boldsymbol{\theta}_t, \boldsymbol{\eta}_t) \quad (63)$$

Such that estimated distribution:

$$\begin{aligned}p(\mathbf{x}_{t+1}^T \boldsymbol{\theta}_{t+1}) &\sim \mathcal{N}(\mathbf{x}_{t+1}^T \boldsymbol{\theta}_t, \mathbf{x}_{t+1}^T \boldsymbol{\eta}_t \mathbf{x}_{t+1}) \\ p(\mathbf{x}_t^T \boldsymbol{\theta}_{t+1}) &\sim \mathcal{N}(\mathbf{x}_t^T \boldsymbol{\theta}_t, \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t)\end{aligned}\quad (64)$$

Given \mathbf{x}_t and \mathbf{x}_{t+1} are drawn from $\mathcal{P}_{\mathcal{X}}$, the k_{t+1} and k_t can be:

$$k_{t+1} = k_t = \mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t \quad (65)$$

To avoid the zero solution and allow manual adjustment, (65) is modified as follows:

$$k_{t+1} = k_t = \kappa \cdot (\mathbf{x}_t^T \boldsymbol{\eta}_t \mathbf{x}_t + \sigma) \quad (66)$$

where σ is a small positive number (e.g. 10^{-5}) and κ is a positive constant (e.g. 1.0).

Proof finished. \square

A.7 OTHER RESULTS OF NUMERICAL SIMULATION

Figures are shown in Figure 7.

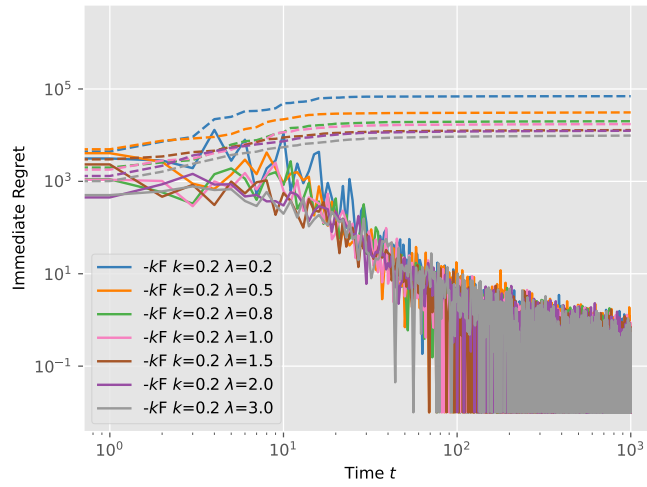
A.8 SUPPLEMENTARY MATERIAL TO CASE 2

Evaluation metrics: We introduced 6 metrics that characterized immediate testset accuracy, incremental task accuracy, knowledge retention ability, degree of knowledge loss, immediate regret, and Kullback-Leibler divergence (KL).

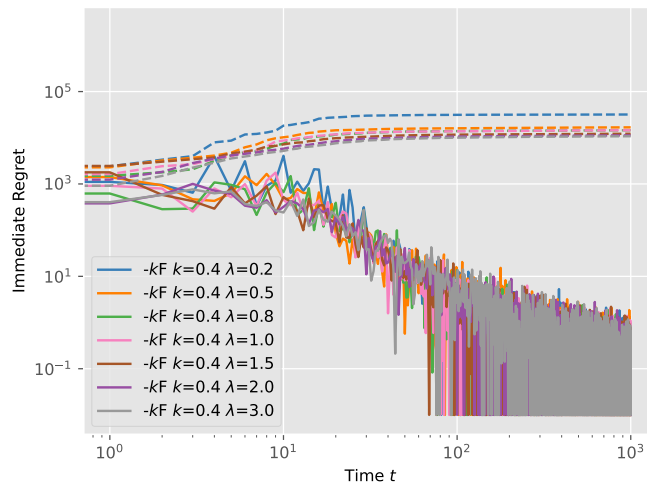
Average task accuracy (ACC) is defined in CL literature as the average accuracy of all previously learned tasks.

$$ACC = \frac{1}{|Q|} \sum_{q=1}^Q R_{Q,q} \quad (67)$$

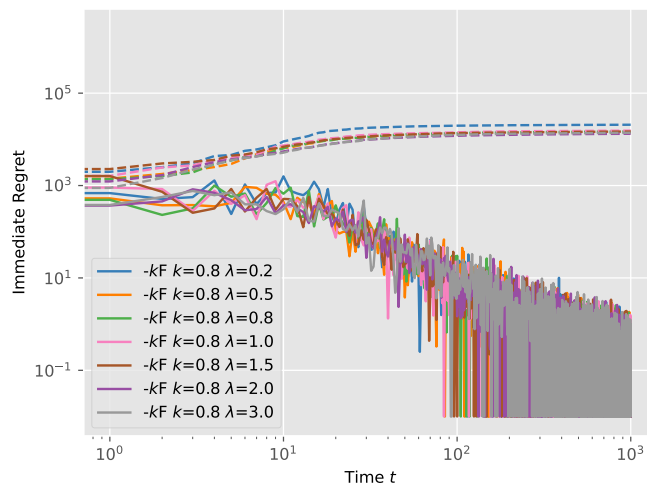
972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025



(a) Relative regrets of online learner using $-0.2F$.



(b) Relative regrets of online learner using $-0.4F$.



(c) Relative regrets of online learner using $-0.8F$.

Table 1: Description of 8 UCI tabular classification datasets

DATASET	$ \mathcal{X} $	$ d $	$ \mathcal{Y} $
page blocks	5473	10	5
glass	214	9	7
image segmentation	2310	9	7
wine quality white	4898	11	7
pendigits	10992	16	10
yeast	1484	8	10
letters	20000	16	26
plant margin	1600	64	100

where $R_{Q,q}$ is the classification accuracy of the learner on task q after learning on task Q ($Q \geq q$). It reflects the task-wise accuracy variation in both CIL and OTCIL processes.

Backward transfer (BWT) is defined as the average difference between the accuracy of all tasks completion and the first learning of one task:

$$BWT = \frac{1}{|Q| - 1} \sum_{q=1}^{Q-1} R_{Q,q} - R_{q,q} \quad (68)$$

BWT indicates the knowledge retention ability of the algorithm where larger values are desired.

Forward transfer (FWT) is defined as the average difference between the accuracy of the first learning of one task and using an independent expert on it:

$$FWT = \frac{1}{|Q| - 1} \sum_{q=2}^Q R_{q,q} - R_q^{ind} \quad (69)$$

where R_q^{ind} denotes the testing accuracy of an independent expert trained only on task q . Higher FWT indicates the learner can acquire more knowledge from newly seen tasks.

Immediate accuracy (acc.(t)) denotes the immediate testing accuracy on entire testset after the learner finishing t -th learning on \mathcal{T} in OTCIL.

$$acc.(t) = R_{t,1..Q} \quad 1 \leq t \leq T \quad (70)$$

It is used to study the dynamic learning performance of using $-kF$ and $-kF$ -Bayes strategies precisely. Note that usually $acc.(T) \leq ACC$ for OTCIL algorithms.

Immediate regret (regret(t)) is defined as the real-time testset loss incurred by the learner:

$$regret(t) = \left\| \frac{\sum_{l=1}^L softmax(\mathbf{X}_{l,te} \boldsymbol{\theta}_{l,t}) - L \cdot \mathbf{Y}_{te}}{L \cdot |\mathcal{X}_{te}|} \right\|_F^2 \quad (71)$$

where $1 \leq t \leq T$ and the subscript still denotes Frobenius norm. Also, we offer cumulative regret to describe the total incurred loss on testset in OTCIL process. Values may be given in logarithm.

Table 2: Testset accuracy comparison of CIL algorithms on 8 UCI tabular datasets

SCENARIO ALGORITHM DATASET	OUR OTCIL LEARNERS			CIL METHODS			OTCIL METHODS		
	edRVFL-R [†] (w/o) <i>acc.(T)%</i> <i>std.%</i>	edRVFL-kF [†] (w/o) <i>acc.(T)%</i> <i>std.%</i>	edRVFL-kF-Bayes [†] (w/o) <i>acc.(T)%</i> <i>std.%</i>	EWC(w/o) <i>ACC%</i> <i>std.%</i>	CRNet-I(w) <i>ACC%</i> <i>std.%</i>	DYSON [†] (w) <i>acc.(T)%</i> <i>std.%</i>			
page blocks	95.68	96.08	95.90	90.13	91.15	91.52			
glass	68.56	68.56	68.86	56.60	66.04	64.15			
image segmentation	89.03	89.03	<u>89.21</u>	61.51	71.02	72.57			
wine quality white	60.29	60.75	<u>61.03</u>	43.61	47.22	56.60			
pendigits	96.33	96.89	<u>97.51</u>	79.40	81.36	78.74			
yeast	59.79	59.97	58.81	45.42	49.73	52.34			
letters	93.27	93.61	94.36	59.81	60.26	82.67			
plant margin	77.92	80.38	79.83	38.75	36.50	53.77			

Note: The † denotes this algorithm can be applied to OTCIL scenario; (w) or (w/o) indicates algorithms with or without using pre-trained models; higher *acc.(T)%* with lower *std.%* is better; the best accuracy of each dataset is shown in **bold**; the underline indicates that the edRVFL-kF-Bayes outperforms both edRVFL-R and edRVFL-kF.