

Supplementary Material: Cost-aware Online Kernel Learning on Imbalanced Data

ACM Reference Format:

. 2021. Supplementary Material: Cost-aware Online Kernel Learning on Imbalanced Data. In *SIGKDD '21: Cost-aware Online Kernel Learning on Imbalanced Data*, August 14–18, 2021, Singapore, Singapore. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/1122445.1122456>

Experimental Results

The experimental results on evaluation of cost with varying budgets are shown in Figure 1. The proposed algorithm Arks achieve a promising result on most of varied setting, and demonstrate the effectiveness and robustness of this proposed algorithm.

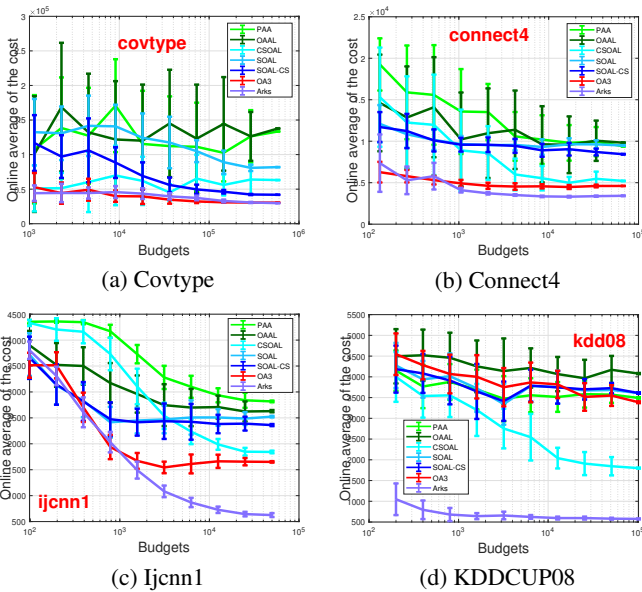


Figure 1: Evaluation of cost with varying budgets.

The performance with Query biases is shown on Figure 2. For Connect4, better performances are obtained when $h_+ \in \{10^4, 10^5\}$ and $h_- \in \{10^{-2}, 10^{-1}, 1\}$ (i.e., upper middle part in Figure 2 (a)). With an imbalanced ratio in Connect4 (i.e., #Pos:#Neg =1:2), the model can boost the perform via biasing to the minority class. For the balanced data Covtype (i.e., #Pos:#Neg =1:1.1), the better performances are usually achieved under symmetric small biases (i.e., bottom left part in Figure 2 (b)). This observation also indicates that

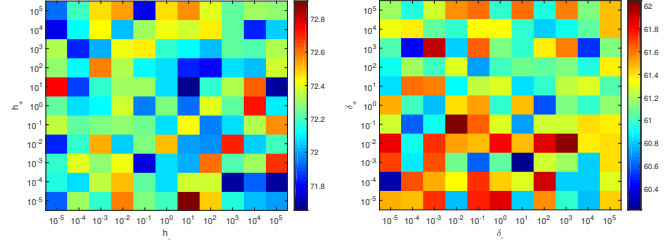
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGKDD '21, August 14–18, 2021, Singapore, Singapore

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/21/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>



(a) Connect4

(b) Covtype

Figure 2: Performance of Sum with varying Query Biases.

this algorithm could perform decently via querying a small amount of labels.

Proof of Theorem

PROPOSITION 1. For any $t > 1$, $\bar{\Phi}_t = [\bar{\phi}_1, \dots, \bar{\phi}_t] \in \mathbb{R}^{D \times t}$, $\bar{\mathbf{K}}_t = \bar{\Phi}_t^\top \bar{\Phi}_t \in \mathbb{R}^{t \times t}$, $\bar{\mathbf{y}}_t = [\bar{y}_1, \dots, \bar{y}_t]^\top$ and $b > 0$. Then the problem (5) could be solved with an optimal solution,

$$\mathbf{w}_t = \bar{\Phi}_{t-1} \left(\bar{\mathbf{K}}_{t-1} + b\mathbf{I} \right)^{-1} \bar{\mathbf{y}}_{t-1}, \quad (1)$$

where the predicted label $\hat{y}_t = \text{sgn}(\mathbf{w}_t^\top \phi_t)$ can be computed without explicit representation ϕ_t ,

$$\hat{y}_t = \text{sgn} \left(\bar{\mathbf{y}}_{t-1}^\top \left(\bar{\mathbf{K}}_{t-1} + b\mathbf{I} \right)^{-1} \bar{\mathbf{k}}_{[t-1],t} \right).$$

PROOF. Given that $\bar{\Phi}_T = [\bar{\phi}_1, \dots, \bar{\phi}_T] \in \mathbb{R}^{D \times T}$ with $\bar{\phi}_t = \sqrt{a_t} \phi_t$ and $\bar{\mathbf{y}}_T = [\bar{y}_1, \dots, \bar{y}_T] \in \mathbb{R}^T$ with $\bar{y}_t = \sqrt{a_t} y_t$, we have

$$\begin{aligned} G_T(\mathbf{w}) &= b\|\mathbf{w}\|^2 + \sum_{t=1}^T a_t (y_t - \mathbf{w}^\top \phi_t)^2 \\ &= b\|\mathbf{w}\|^2 + \sum_{t=1}^T a_t \left(y_t^2 - 2(y_t \mathbf{w}^\top \phi_t) + \mathbf{w}^\top \phi_t \phi_t^\top \mathbf{w} \right) \\ &= \mathbf{w}^\top \left(b\mathbf{I} + \sum_{t=1}^T a_t \phi_t \phi_t^\top \right) \mathbf{w} - 2\mathbf{w}^\top \left(\sum_{t=1}^T a_t \phi_t y_t \right) + \sum_{t=1}^T a_t y_t^2 \\ &= \mathbf{w}^\top (b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top) \mathbf{w} - 2\mathbf{w}^\top (\bar{\Phi}_T \bar{\mathbf{y}}_T) + \bar{\mathbf{y}}_T^\top \bar{\mathbf{y}}_T. \end{aligned}$$

Then follows that $\nabla G_T(\mathbf{w}) = 2(b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top) \mathbf{w} - 2\bar{\Phi}_T \bar{\mathbf{y}}_T$, $\nabla^2 G_T(\mathbf{w}) = 2(b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top) > 0$. Thus $G_T(\mathbf{w})$ is convex and it is minimal if $\nabla G_T(\mathbf{w}) = (b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top) \mathbf{w} - \bar{\Phi}_T \bar{\mathbf{y}}_T = 0$ with

$$\mathbf{w} = (b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top)^{-1} \bar{\Phi}_T \bar{\mathbf{y}}_T.$$

For any $\bar{\Phi}_T = [\bar{\phi}_1, \dots, \bar{\phi}_T] \in \mathbb{R}^{D \times T}$ matrix and $b > 1$,

$$\bar{\Phi}_T \bar{\Phi}_T^\top (\bar{\Phi}_T \bar{\Phi}_T^\top + b\mathbf{I}_D)^{-1} = \bar{\Phi}_T (\bar{\Phi}_T^\top \bar{\Phi}_T + b\mathbf{I}_T)^{-1} \bar{\Phi}_T^\top = \bar{\Phi}_T (\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} \bar{\Phi}_T^\top$$

With above equation, we have

$$\begin{aligned} (\bar{\Phi}_T \bar{\Phi}_T^\top + b\mathbf{I}_D)^{-1} &= \frac{1}{b} b\mathbf{I}_D (\bar{\Phi}_T \bar{\Phi}_T^\top + b\mathbf{I}_D)^{-1} \\ &= \frac{1}{b} \left(\bar{\Phi}_T \bar{\Phi}_T^\top + b\mathbf{I}_D - \bar{\Phi}_T \bar{\Phi}_T^\top \right) (\bar{\Phi}_T \bar{\Phi}_T^\top + b\mathbf{I}_D)^{-1} \\ &= \frac{1}{b} \left(\mathbf{I}_D - \bar{\Phi}_T \bar{\Phi}_T^\top (\bar{\Phi}_T \bar{\Phi}_T^\top + b\mathbf{I}_D)^{-1} \right) = \frac{1}{b} \left(\mathbf{I}_D - \bar{\Phi}_T (\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} \bar{\Phi}_T^\top \right) \end{aligned}$$

Substituting the two conclusions, we obtain

$$\begin{aligned} \mathbf{w}_{T+1} &= (b\mathbf{I}_D + \bar{\Phi}_T \bar{\Phi}_T^\top)^{-1} \bar{\Phi}_T \bar{\mathbf{y}}_T \\ &= \frac{1}{b} \left(\mathbf{I}_D - \bar{\Phi}_T (\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} \bar{\Phi}_T^\top \right) \bar{\Phi}_T \bar{\mathbf{y}}_T \\ &= \frac{1}{b} \bar{\Phi}_T \left(\bar{\mathbf{y}}_T - (\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} \bar{\mathbf{K}}_T \bar{\mathbf{y}}_T \right) \\ &= \frac{1}{b} \bar{\Phi}_T \left(\mathbf{I}_T - (\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} \bar{\mathbf{K}}_T \right) \bar{\mathbf{y}}_T \\ &= \frac{1}{b} \bar{\Phi}_T \left((\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} (\bar{\mathbf{K}}_T + b\mathbf{I}_T - \bar{\mathbf{K}}_T) \right) \bar{\mathbf{y}}_T \\ &= \bar{\Phi}_T (\bar{\mathbf{K}}_T + b\mathbf{I}_T)^{-1} \bar{\mathbf{y}}_T. \end{aligned}$$

Substituting the \mathbf{w}_T back to $\hat{\mathbf{y}}_T = \text{sgn}(\mathbf{w}_T^\top \phi_T)$, we finish the proof. \square

Proof on Lemma 1

LEMMA 1. Let $(\phi_1, y_1), \dots, (\phi_T, y_T)$ be a sequence of input samples, where $\phi_t \in \mathcal{H}$ and $y_t \in \{\pm 1\}$ for all t . The model parameter \mathbf{w}_t is learned by Eq. (1) with $b > 0$. For any $\mathbf{w} \in \mathcal{H}$, it satisfies:

$$\begin{aligned} \sum_{t=1}^T a_t (y_t - \mathbf{w}_t^\top \phi_t)^2 - \sum_{t=1}^T a_t (y_t - \mathbf{w}^\top \phi_t)^2 \\ \leq b \|\mathbf{w}\|^2 + 2\alpha(1 + C^2) \log \det\left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I}\right), \end{aligned} \quad (2)$$

where $\alpha = \frac{\mu_p T_n}{\mu_n T_p}$ or $\frac{c_p}{c_n}$ and $|\mathbf{w}_t^\top \phi_t| \leq C$ for $t \in [T]$.

PROOF. Defined that $\mathbf{A}_T = (b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top)$ with $\mathbf{A}_T = \mathbf{A}_{T-1} + \bar{\phi}_T \bar{\phi}_T^\top$, and $\mathbf{b}_T = \bar{\Phi}_T \bar{\mathbf{y}}_T$ with $\mathbf{b}_T = \mathbf{b}_{T-1} + \bar{y}_T \bar{\phi}_T$, so that $\mathbf{w}_T = \mathbf{A}_{T-1}^{-1} \mathbf{b}_{T-1}$. According to the Woodbury formulation, $\mathbf{A}_T^{-1} = \mathbf{A}_{T-1}^{-1} - \frac{\mathbf{A}_{T-1}^{-1} \bar{\phi}_T \bar{\phi}_T^\top \mathbf{A}_{T-1}^{-1}}{1 + \bar{\phi}_T^\top \mathbf{A}_{T-1}^{-1} \bar{\phi}_T}$, and with annotation $\tau_t = \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1} \bar{\phi}_t$, we have

$$\bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t = \bar{\phi}_t^\top \left(\mathbf{A}_{t-1}^{-1} - \frac{\mathbf{A}_{t-1}^{-1} \bar{\phi}_t \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1}}{1 + \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1} \bar{\phi}_t} \right) \bar{\phi}_t = \frac{\tau_t}{1 + \tau_t}.$$

Substituting $\mathbf{w}_T = \mathbf{A}_{T-1}^{-1} \mathbf{b}_{T-1}$ back to F_T , $F_T = -\mathbf{b}_T^\top \mathbf{A}_{T-1}^{-1} \mathbf{b}_T + \bar{\mathbf{y}}_T^\top \bar{\mathbf{y}}_T$. It derives that

$$\begin{aligned} F_{T-1} - F_T &= \left(-\mathbf{b}_{T-1}^\top \mathbf{A}_{T-1}^{-1} \mathbf{b}_{T-1} + \bar{\mathbf{y}}_{T-1}^\top \bar{\mathbf{y}}_{T-1} \right) - \left(-\mathbf{b}_T^\top \mathbf{A}_T^{-1} \mathbf{b}_T + \bar{\mathbf{y}}_T^\top \bar{\mathbf{y}}_T \right) \\ &= (\mathbf{b}_{T-1} + \bar{y}_T \bar{\phi}_T)^\top \mathbf{A}_T^{-1} (\mathbf{b}_{T-1} + \bar{y}_T \bar{\phi}_T) - \mathbf{b}_{T-1}^\top \mathbf{A}_{T-1}^{-1} \mathbf{b}_{T-1} - \bar{y}_T^2 \\ &= \bar{y}_T^2 \left(\bar{\phi}_T^\top \mathbf{A}_T^{-1} \bar{\phi}_T - 1 \right) + 2\bar{y}_T \mathbf{b}_{T-1}^\top \mathbf{A}_T^{-1} \bar{\phi}_T + \mathbf{b}_{T-1}^\top (\mathbf{A}_T^{-1} - \mathbf{A}_{T-1}^{-1}) \mathbf{b}_{T-1} \\ &= \frac{-a_T y_T^2}{1 + \tau_T} + \frac{2a_T y_T f_T}{1 + \tau_T} + \frac{-a_T f_T^2}{1 + \tau_T} \\ &= -a_T (y_T - f_T)^2 (1 - \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t) \end{aligned}$$

where the fourth equation holds since $f_t = \mathbf{w}_t^\top \phi_t = \mathbf{b}_t^\top \mathbf{A}_t^{-1} \phi_t$, and $\bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t = \frac{\tau_t}{1 + \tau_t}$, and the fifth equation holds since $\frac{1}{1 + \tau_t} = 1 - \frac{\tau_t}{1 + \tau_t} = 1 - \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t$. Summing over t from 1 to T , it yields $\sum_{t=1}^T (F_t - F_{t-1}) = F_T = \min_{\mathbf{w} \in \mathcal{H}} \left(b \|\mathbf{w}\|^2 + \sum_{t=1}^T a_t (y_t - \mathbf{w}^\top \phi_t)^2 \right)$. For any $\mathbf{w} \in \mathcal{H}$, we have

$$\sum_{t=1}^T a_t (y_t - f_t)^2 (1 - \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t) \leq b \|\mathbf{w}\|^2 + \sum_{t=1}^T a_t (y_t - \mathbf{w}^\top \phi_t)^2$$

Rearrange the terms, and bound $|\mathbf{w}_t^\top \phi_t| \leq C$ and $a_t \leq \alpha$ for $t \in [T]$,

$$\begin{aligned} \sum_{t=1}^T a_t (y_t - \mathbf{w}_t^\top \phi_t)^2 - a_t (y_t - \mathbf{w}^\top \phi_t)^2 \\ \leq b \|\mathbf{w}\|^2 + 2\alpha(1 + C^2) \sum_{t=1}^T \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t \end{aligned} \quad (3)$$

Since $\mathbf{A}_t = \mathbf{A}_{t-1} + \bar{\phi}_t \bar{\phi}_t^\top \Rightarrow \mathbf{A}_{t-1} \mathbf{A}_t^{-1} = \mathbf{A}_t \mathbf{A}_{t-1}^{-1} - \bar{\phi}_t \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1}$, according to Sylvester's determinant theorem,

$$\det(\mathbf{A}_t \mathbf{A}_{t-1}^{-1} - \bar{\phi}_t \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1}) = \det(\mathbf{I}_D - \bar{\phi}_t \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1}) = 1 - \bar{\phi}_t^\top \mathbf{A}_{t-1}^{-1} \bar{\phi}_t,$$

while $\det(\mathbf{A}_{t-1} \mathbf{A}_t^{-1}) = \det(\mathbf{A}_{t-1}) \det(\mathbf{A}_t)^{-1}$. With these equations, we obtain $\bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t = 1 - \frac{\det(\mathbf{A}_{t-1})}{\det(\mathbf{A}_t)} \leq \log \left(\frac{\det(\mathbf{A}_t)}{\det(\mathbf{A}_{t-1})} \right)$, where the inequality holds since $1 - \frac{1}{x} \leq \log x$ for $x > 0$. Summing over $t = 1, \dots, T$,

$$\begin{aligned} \sum_{t=1}^T \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t &\leq \sum_{t=1}^T \log \left(\frac{\det(\mathbf{A}_t)}{\det(\mathbf{A}_{t-1})} \right) = \log \left(\frac{\det(\mathbf{A}_T)}{\det(\mathbf{A}_0)} \right) \\ &= \log \left(\det \left(\frac{1}{b} \bar{\Phi}_T \bar{\Phi}_T^\top + \mathbf{I}_D \right) \right) = \log \det \left(\frac{1}{b} \bar{\Phi}_T \bar{\Phi}_T^\top + \mathbf{I}_T \right). \end{aligned}$$

where the last equality holds according to Sylvester's determinant theorem, and we obtain $\sum_{t=1}^T \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t \leq \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right) \right)$. Substituting this bound back to Eq. (3), it concludes the proof. \square

Proof of Theorem 1

THEOREM 1. Given an arbitrary sequence $\{(\phi_t, y_t)\}_{t=1}^T$, the algorithm learns on only queried trails $\{Z_t \phi_t\}_{t=1}^T$ where $Z_t \sim \mathcal{B}(1, p_t)$ with $p_t = \frac{h_+}{h_+ + \max(0, \Theta_t)}$ if $\hat{f}_t \geq 0$ and $p_t = \frac{h_-}{h_- + \max(0, \Theta_t)}$ if $\hat{f}_t < 0$. Let $\ell_h(\cdot)$ be the hinge loss, set $h_+ = \sqrt{\frac{\alpha \log \det(\frac{1}{b} \bar{\mathbf{K}}_{T_p} + \mathbf{I})}{(b + T_p) C^2}}$ and $h_- = \sqrt{\frac{\alpha \log \det(\frac{1}{b} \bar{\mathbf{K}}_{T_n} + \mathbf{I})}{(b + T_n) C^2}}$. For any $\mathbf{w} \in \mathcal{H}$, the expected weighted mistake number is bounded by:

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T a_t M_t \right] &= \mathbb{E} \left[\sum_{t=1, y_t=+1}^T a_t M_t + \sum_{t=1, y_t=-1}^T a_t M_t \right] \\ &\leq \sum_{t=1}^T a_t \ell_h(\mathbf{w}) + C \sqrt{\alpha(b + T) \log \det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right)} \end{aligned}$$

PROOF. According to the Lemma 1, for any $\mathbf{w} \in \mathcal{H}$,

$$\sum_{t=1}^T \frac{a_t (y_t - f_t)^2}{1 + \tau_t} \leq \mathbf{w}^\top (b\mathbf{I} + \bar{\Phi}_T \bar{\Phi}_T^\top) \mathbf{w} + \sum_{t=1}^T a_t (y_t^2 - 2y_t \mathbf{w}^\top \phi_t)$$

Since \mathbf{w} is a random variable, we use $h\mathbf{w}$ to replace \mathbf{w} , and add $\sum_{t=1}^T 2a_t h$ on both hands of the inequality,

$$\begin{aligned} & \sum_{t=1}^T 2a_t \left(h + \frac{-y_t f_t - \tau_t/2}{1 + \tau_t} \right) \\ & \leq h^2 \mathbf{w}^\top \mathbf{A}_T \mathbf{w} + \sum_{t=1}^T 2h a_t (1 - y_t \mathbf{w}^\top \phi_t) - \sum_{t=1}^T \frac{a_t f_t^2}{1 + \tau_t} \end{aligned}$$

Since $1 - y_t \mathbf{w}^\top \phi_t \leq \max(0, 1 - y_t \mathbf{w}^\top \phi_t) = \ell_t(\mathbf{w})$, where ℓ_t is the hinge loss, we obtain

$$\sum_{t=1}^T a_t \left(h + \frac{-y_t f_t - \tau_t/2}{1 + \tau_t} \right) \leq \frac{h^2}{2} \mathbf{w}^\top \mathbf{A}_T \mathbf{w} + h \sum_{t=1}^T a_t \ell_t(\mathbf{w}) - \sum_{t=1}^T \frac{a_t f_t^2}{2(1 + \tau_t)}$$

Omitting the term $-\sum_{t=1}^T \frac{a_t f_t^2}{2(1 + \tau_t)}$ that does not affect inequality, and considering the inequality with the scenarios $M_t Z_t$,

$$\sum_{t=1}^T M_t Z_t a_t \left(h + \frac{-y_t f_t - \tau_t/2}{1 + \tau_t} \right) \leq \frac{h^2}{2} \mathbf{w}^\top \mathbf{A}_T \mathbf{w} + h \sum_{t=1}^T a_t \ell_t(\mathbf{w}),$$

One can easily prove that this inequality still holds for $M_t Z_t = 0$. If an error occurs at trial t , i.e., $M_t Z_t = 1$, we have $-y_t f_t = |f_t|$. Then the left hand of the inequality becomes

$$\sum_{t=1}^T M_t Z_t a_t \left(\frac{|f_t| - \tau_t/2}{1 + \tau_t} \right) = \sum_{t=1}^T M_t Z_t a_t \Theta_t$$

Considering that the algorithm queries an input and suffers a mistake at round t , i.e., $Z_t = 1$ and $M_t = 1$. There are two mistake cases: (1) False negative: true label $y_t = +1$ while the prediction $f_t < 0$; (2) False positive: true label $y_t = -1$ while the prediction $f_t \geq 0$. For the first case, we have:

$$\begin{aligned} & \sum_{t, y_t=+1} M_t Z_t a_t (h_+ + \Theta_t) \\ & \leq \sum_{t, y_t=+1} h_+ a_t \ell_h(\mathbf{w}) + \frac{h_+^2}{2} \mathbf{w}^\top \left(b\mathbf{I} + \sum_{t, y_t=+1} \bar{\phi}_t \bar{\phi}_t^\top \right) \mathbf{w} \end{aligned}$$

Now we would like to remove the random variable Z_t . First, when the confidence score $\Theta_t > 0$, taking the expectation over variables $\mathbb{E}(Z_t) = \frac{h_+}{h_+ + \Theta_t}$, we have:

$$\mathbb{E} \left[\sum_{t, y_t=+1} h_+ M_t a_t \right] \leq \frac{h_+^2}{2} \mathbf{w}^\top \left(b\mathbf{I} + \sum_{t, y_t=+1} \bar{\phi}_t \bar{\phi}_t^\top \right) \mathbf{w} + \sum_{t=1, y_t=+1}^T h_+ a_t \ell_h(\mathbf{w})$$

In addition, one can easily prove this inequality holds for $M_t = 0$. Second, when $\Theta_t \leq 0$, the random variable is assigned to $\mathbb{E}(Z_t) = 1$ and $M_t \in \{0, 1\}$,

$$\begin{aligned} & \mathbb{E}[M_t Z_t a_t (h_+ + \Theta_t)] = \mathbb{E}[Z_t] \mathbb{E} \left[M_t a_t \left(h_+ + \frac{|f_t| - \tau_t/2}{1 + \tau_t} \right) \right] \\ & \geq \mathbb{E} \left[M_t a_t \left(h_+ - \frac{\tau_t/2}{1 + \tau_t} \right) \right] \geq \mathbb{E}[h_+ M_t a_t] - \mathbb{E} \left[\frac{a_t \tau_t}{2(1 + \tau_t)} \right] \\ & = \mathbb{E}[h_+ M_t a_t] - \frac{a_t}{2} \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t. \end{aligned}$$

Combining above two scenarios for Θ_t , we obtain

$$\begin{aligned} \mathbb{E} \left[\sum_{t, y_t=+1} a_t M_t \right] & \leq \frac{h_+}{2} \mathbf{w}^\top \left(b\mathbf{I} + \sum_{t, y_t=+1} \bar{\phi}_t \bar{\phi}_t^\top \right) \mathbf{w} + \sum_{t, y_t=+1} a_t \ell_h(\mathbf{w}) \\ & \quad + \sum_{t, y_t=+1, \Theta_t \leq 0} \frac{a_t}{2h_+} \bar{\phi}_t^\top \mathbf{A}_t^{-1} \bar{\phi}_t \\ & \leq \frac{h_+}{2} (b + T_p) C^2 + \sum_{t, y_t=+1} a_t \ell_h(\mathbf{w}) + \frac{\alpha}{2h_+} \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_{T_p} + \mathbf{I} \right) \right) \end{aligned}$$

where $\|\mathbf{w}\|_2 \leq C$ and T_p is number of positive labels in the last

inequality. If we set $h_+ = \sqrt{\frac{\alpha \log(\det(\frac{1}{b} \bar{\mathbf{K}}_{T_p} + \mathbf{I}))}{(b + T_p) C^2}}$, then

$$\mathbb{E} \left[\sum_{t, y_t=+1} a_t M_t \right] \leq \sum_{t, y_t=+1} a_t \ell_h(\mathbf{w}) + C \sqrt{\alpha (b + T_p) \log \det \left(\frac{1}{b} \bar{\mathbf{K}}_{T_p} + \mathbf{I} \right)} \quad (4)$$

When $y_t = -1$, setting $h_- = \sqrt{\frac{\alpha \log(\det(\frac{1}{b} \bar{\mathbf{K}}_{T_n} + \mathbf{I}))}{(b + T_n) C^2}}$, we have

$$\mathbb{E} \left[\sum_{t, y_t=-1} a_t M_t \right] \leq \sum_{t, y_t=-1} a_t \ell_h(\mathbf{w}) + C \sqrt{\alpha (b + T_n) \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_{T_n} + \mathbf{I} \right) \right)} \quad (5)$$

Summing Eq. (4) and (5) will give:

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T a_t M_t \right] & = \mathbb{E} \left[\sum_{t=1, y_t=+1}^T a_t M_t + \sum_{t=1, y_t=-1}^T a_t M_t \right] \\ & \leq \sum_{t=1}^T a_t \ell_h(\mathbf{w}) + 2C \sqrt{\alpha (b + T) \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right) \right)} \end{aligned}$$

Then, we conclude the proofs of Theorem. \square

Proof of Theorem 2

THEOREM 2. Under the same condition in Theorem 1, by setting $\alpha = \frac{\mu_p T_n}{\mu_n T_p}$, the proposed algorithm satisfies for any $\mathbf{w} \in \mathcal{H}$:

$$\mathbb{E}[\text{sum}] \geq 1 - \frac{\mu_n}{T_n} \left[\sum_{t=1}^T a_t \ell_h(\mathbf{w}) + 2C \sqrt{\frac{\mu_p T_n}{\mu_n T_p} (b + T) \log \det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right)} \right]$$

By setting $\mu_p = \mu_n = 0.5$, we can easily obtain the bound of the balanced accuracy.

PROOF. Associating the cost-aware sum with the cost-aware loss function, we have

$$\begin{aligned} \text{sum} & = \mu_p \frac{T_p - M_p}{T_p} + \mu_n \frac{T_n - M_n}{T_n} \\ & = 1 - \frac{\mu_n}{T_n} \left[\sum_{y_i=+1} \frac{\mu_p T_n}{\mu_n T_p} \mathbb{I}(\hat{y}_i < 0) + \sum_{y_i=-1} \mathbb{I}(\hat{y}_i \geq 0) \right]. \end{aligned}$$

By setting $a_t = \frac{\mu_p T_n}{\mu_n T_p} \mathbb{I}(y_t=+1) + \mathbb{I}(y_t=-1)$, we have

$$\begin{aligned} \mathbb{E}[\text{sum}] &= 1 - \frac{\mu_n}{T_n} \mathbb{E} \left[\sum_{t=1}^T a_t M_t \right] \\ &\geq 1 - \frac{\mu_n}{T_n} \left[\sum_{t=1}^T a_t \ell_h(\mathbf{w}) + 2C \sqrt{\frac{\mu_p T_n}{\mu_n T_p} (b+T) \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right) \right)} \right] \end{aligned}$$

□

Proof of Theorem 3

THEOREM 3. *Under the same condition in Theorem 1, by setting $\alpha = \frac{c_p}{c_n}$, the proposed algorithm satisfies for any $\mathbf{w} \in \mathcal{H}$:*

$$\mathbb{E}[\text{cost}] \leq c_n \left[\sum_{t=1}^T a_t \ell_h(\mathbf{w}) + 2C \sqrt{\frac{c_p}{c_n} (b+T) \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right) \right)} \right]$$

By setting $c_p = c_n$, we can easily obtain the bound of the balanced penalty.

PROOF. Associating the cost-aware cost with the cost-aware loss function, we have

$$\text{cost} = c_p \times M_p + c_n \times M_n = c_n \left[\sum_{y_i=+1} \frac{c_p}{c_n} \mathbb{I}(\hat{y}_i < 0) + \sum_{y_i=-1} \mathbb{I}(\hat{y}_i \geq 0) \right]$$

By setting $a_t = \frac{c_p}{c_n} \mathbb{I}(y_t=+1) + \mathbb{I}(y_t=-1)$, we have

$$\begin{aligned} \mathbb{E}[\text{cost}] &= c_n \mathbb{E} \left[\sum_{t=1}^T a_t M_t \right] \\ &\leq c_n \left[\sum_{t=1}^T a_t \ell_h(\mathbf{w}) + 2C \sqrt{\frac{c_p}{c_n} (b+T) \log \left(\det \left(\frac{1}{b} \bar{\mathbf{K}}_T + \mathbf{I} \right) \right)} \right] \end{aligned}$$

□