

Online Pairwise Learning Algorithms with Convex Loss Functions[☆]

Junhong Lin, Yunwen Lei*, Bo Zhang*, Ding-Xuan Zhou*

Department of Mathematics, City University of Hong Kong, Kowloon, Hong Kong, China

Abstract

Online pairwise learning algorithms with general convex loss functions without regularization in a Reproducing Kernel Hilbert Space (RKHS) are investigated. Under mild conditions on loss functions and the RKHS, upper bounds for the expected excess generalization error are derived in terms of the approximation error when the stepsize sequence decays polynomially. In particular, for Lipschitz loss functions such as the hinge loss, the logistic loss and the absolute-value loss, the bounds can be of order $O(T^{-\frac{1}{3}} \log T)$ after T iterations, while for the least squares loss, the bounds can be of order $O(T^{-\frac{1}{4}} \log T)$. In comparison with previous works for these algorithms, a broader family of convex loss functions is studied here, and refined upper bounds are obtained.

Keywords:

Learning theory, Online learning, Learning Theory, Reproducing Kernel Hilbert Space, Pairwise learning

1. Introduction

Many classical learning tasks can be modeled as learning a good estimator or predictor $f : X \rightarrow Y$ based on an observed dataset $\{(x_t, y_t)\}_{t=1}^T$ of input-output

[☆]The work described in this paper is supported partially by the Research Grants Council of Hong Kong [Project No. CityU 104113]. The corresponding author is Yunwen Lei. Junhong Lin is now within the LCSL, MIT & Istituto Italiano di Tecnologia, Cambridge, MA 02139, USA

*Corresponding author

Email addresses: jhlin5@hotmail.com (Junhong Lin), yunweilei@cityu.edu.hk (Yunwen Lei), bozhang37-c@my.cityu.edu.hk (Bo Zhang), mazhou@cityu.edu.hk (Ding-Xuan Zhou)

4 samples from $X \times Y$, where X is an input space and $Y \subseteq \mathbb{R}$ an output space.
 5 Learning algorithms are often implemented by minimizing $\frac{1}{T} \sum_{t=1}^T V(y_t, f(x_t))$
 6 over a hypothesis space of functions in various ways including regularization
 7 schemes [26]. Here $V : \mathbb{R}^2 \rightarrow \mathbb{R}_+$ is a *loss function* used for measuring the perfor-
 8 mance of a predictor f . It induces a *local error* $V(y, f(x))$ over an input-output
 9 sample $(x, y) \in X \times Y$. For non-parametric regression with $Y = \mathbb{R}$, the least
 10 squares loss function $V(y, a) = (y - a)^2$ is often used and, for an input $x \in X$ and
 11 an estimator f , the induced local error $V(y, f(x)) = (y - f(x))^2$ measures how
 12 well the predicted value $f(x)$ approximates the output value $y \in \mathbb{R}$. For binary
 13 classification with $Y = \{1, -1\}$ consisting of the two labels corresponding to the
 14 two classes, the misclassification loss function $V(y, a) = \chi_{(-\infty, 0)}(ya)$ generated
 15 by the characteristic function of the interval $(-\infty, 0)$ is a natural choice, and the
 16 induced local error $V(y, f(x)) = \chi_{(-\infty, 0)}(yf(x))$ over a sample $(x, y) \in X \times Y$
 17 equals 1 when the sign of $f(x)$ and y correspond to the two different labels in
 18 Y (that is, $yf(x) < 0$), while $V(y, f(x)) = 0$ when they correspond to a same
 19 label with $yf(x) \geq 0$. But the characteristic function $\chi_{(-\infty, 0)}$ is not convex, and
 20 the optimization problems involved in the related learning algorithms are not
 21 convex. For designing efficient learning algorithms, $\chi_{(-\infty, 0)}$ may be replaced
 22 by a convex function $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$, leading to convex optimization problems in-
 23 volving the local error $V(y, f(x)) = \phi(yf(x))$. One choice of ϕ is the hinge loss
 24 $\phi_h(v) = \max\{1 - v, 0\}$ used in the classical support vector machines for solv-
 25 ing binary classification problems [26]. The above learning framework has been
 26 well developed within the last two decades [26, 9]. It might be categorized as
 27 “pointwise learning”, as the local error $V(y, f(x))$ takes only one sample point
 28 $(x, y) \in X \times Y$ into account.

29 In this paper, we study another important family of learning problems cate-
 30 gorized as “pairwise learning” in which the local error takes a pair $\{(x, y), (x', y')\}$
 31 of two samples from $X \times Y$ into account. Its learning tasks include ranking [1, 8],
 32 similarity and metric learning [5, 28], AUC maximization [34], and gradient
 33 learning [20, 30, 19]. The goal of *pairwise learning* is to learn a good predictor
 34 $f : X^2 \rightarrow \mathbb{R}$ predicting a value $f(x, x') \in \mathbb{R}$ for each input pair $(x, x') \in X^2$ ac-

35 cording to various tasks. To measure the learning performance of a predictor f ,
 36 we use a loss function $V : \mathbb{R}^2 \rightarrow \mathbb{R}_+$ to induce the local error $V(r(y, y'), f(x, x'))$
 37 over two input-output samples $(x, y), (x', y') \in X \times Y$, where $r : Y \times Y \rightarrow \mathbb{R}$ is
 38 a function, called *reducing function*, chosen according to the learning task. The
 39 reducing function r is an essential concept making pairwise learning different
 40 from pointwise learning. We demonstrate how to choose the reducing function
 41 r by the following examples.

- 42 1. For the least squares regression with $Y = \mathbb{R}$ and $V(y, a) = (y - a)^2$, a
 43 sample (x, y) is drawn from a probability measure and the expected value
 44 of $y \in \mathbb{R}$ given $x \in X$ equals $f^*(x)$, the value of the conditional mean
 45 (regression) function f^* at x . So $y - y' = f^*(x) - f^*(x')$ in expectation
 46 and we choose the reducing function $r : Y \times Y \rightarrow \mathbb{R}$ as the output value
 47 difference $r(y, y') = y - y'$. Then the local error $V(r(y, y'), f(x, x')) =$
 48 $(y - y' - f(x, x'))^2$ measures how well the predicted value $f(x, x')$ for an
 49 input pair (x, x') approximates $f^*(x) - f^*(x')$ via the output value differ-
 50 ence $y - y'$.
- 51 2. For metric learning in binary classification with $Y = \{1, -1\}$, we aim to
 52 learn a metric f such that a pair (x, x') of inputs (objects) from the same
 53 class ($y = y'$) are close to each other while a pair from different classes ($y \neq$
 54 y') have a large distance $f(x, x')$. A typical choice of the reducing function
 55 $r : Y \times Y \rightarrow \mathbb{R}$ is given by $r(y, y') = 1$ if $y = y'$ and -1 otherwise [5]. The
 56 local error induced by the convex loss function $V(y, a) = \max\{0, 1 + ya\}$
 57 is $V(r(y, y'), f(x, x')) = \max\{0, 1 + r(y, y')f(x, x')\}$. It gives a large local
 58 error $1 + f(x, x')$ if the distance $f(x, x')$ between the input pair (x, x')
 59 from the same class ($y = y'$) is large.
- 60 3. For ranking in a regression framework with $Y = \mathbb{R}$, we aim to learn a good
 61 ordering f between objects (inputs) based on their observed features such
 62 that $f(x, x') < 0$ if x is preferred over x' meaning that the ranking labels
 63 satisfy $y < y'$. A typical choice [21] of the reducing function $r : Y \times Y \rightarrow \mathbb{R}$
 64 is given by $r(y, y') = \text{sign}(y - y')$, the sign of $y - y'$. Then the local

65 error induced by the hinge loss ϕ_h is $V(r(y, y'), f(x, x')) = \phi_h(\text{sign}(y -$
66 $y')f(x, x'))$.

67 Batch learning and online learning are two kinds of learning algorithms. The
68 former uses an entire dataset to perform learning tasks, while the latter uses
69 the dataset in a stream way. For batch learning algorithms in the pairwise
70 learning framework, theoretical error and robustness analysis has been carried
71 out in [1, 8, 21, 5, 7]. One challenge in conducting analysis in pairwise learning
72 is that pairs of training samples are not independent. For example, given the
73 independently and identically distributed (i.i.d.) samples $\{z_t = (x_t, y_t)\}_{t=1}^T$, a
74 batch algorithm for pairwise learning possibly involves a target function

$$\frac{T(T-1)}{2} \sum_{1 \leq i < j \leq T} V(r(y_i, y_j), f(x_i, x_j)) + \text{pen}(f, \lambda), \quad (1.1)$$

75 where $\text{pen}(f, \lambda) \geq 0$ is some regularization term used to avoid overfitting. In this
76 case, local errors $V(r(y_i, y_j), f(x_i, x_j))$ and $V(r(y_i, y_{j'}), f(x_i, x_{j'}))$ are indeed
77 dependent. Thus, standard techniques for classification and regression cannot
78 be directly applied, and new tools such as U-statistics [8] or algorithmic stability
79 [1] are necessary for the analysis.

80 In spite of their good theoretical guarantees, batch algorithms for pairwise
81 learning may be difficult to implement for large-scale learning problems in prac-
82 tice. Indeed, even for the simpler case of pointwise learning, the computational
83 complexity of batch algorithms with many loss functions is of order $O(T^3)$.
84 Moreover, batch algorithms for pairwise learning suffer from extra computa-
85 tional burden of optimizing an objective defined over $O(T^2)$ possible sample
86 pairs.

87 In practical applications, online learning may be more favorable, due to its
88 scalability to large datasets and applicability to situations where the samples
89 are collected sequentially. Theoretical results for online learning in classification
90 and regression have been well developed (see for example [6, 24, 31, 2, 22, 18]
91 and references therein), but there is relatively little work for online learning in
92 pairwise learning. Recent research of this direction can be found in [15, 27,

93 32]. In particular, online pairwise learning in a linear space was investigated in
94 [15, 27], and convergence results were established for the average of the iterates
95 under the assumption of uniform boundedness of the loss function, with a rate
96 $O(1/\sqrt{T})$ in the general convex case, or a rate $O(1/T)$ in the strongly convex
97 case. Online pairwise learning in a RKHS with the least squares loss was studied
98 in [32] where bounds in probability were derived for the excess generalization
99 error.

100 In this paper, we improve the analysis of online pairwise learning (see Al-
101 gorithm 1 in the next section) in a RKHS with general convex loss functions.
102 Our main purpose is to develop convergence results for such learning algorithms
103 using polynomially decaying stepsize sequences. Unlike [15, 27], we do not as-
104 sume that the iterates are restricted to a bounded domain or the loss function is
105 strongly convex. In particular, we will provide bounds for the expected excess
106 generalization error, under a mild condition on approximation errors and an
107 increment condition on the loss. For Lipschitz loss functions such as the hinge
108 loss and the logistic loss, our bounds can be of order $O(T^{-\frac{1}{3}} \log T)$, while for the
109 least squares loss, our bounds can be of order $O(T^{-\frac{1}{4}} \log T)$. For general convex
110 loss functions, previous error analysis techniques dealing with the least squares
111 loss in [32], which rely on integral operators, do not apply and are replaced
112 by tools from convex analysis and Rademacher complexity. The key to our
113 proof is an error decomposition, which enables us to study the weighted excess
114 generalization error in terms of the weighted average and the moving weighted
115 average. The novelty lies in an estimate of the differences between partial and
116 generalization errors of the learning sequence. We shall establish bounds for the
117 learning sequence using tools from convex analysis, and give uniform bounds
118 for the differences between partial and full generalization errors over any given
119 ball using Rademacher complexity. Our methods may be applied to pairwise
120 learning with non-convex loss functions. In particular, it would be interesting
121 to extend our methods to online learning or gradient descent methods for a
122 minimum error entropy principle [10, 14].

123 **2. Main Results with Discussions**

124 In this section, after stating our pairwise learning problems and basic as-
 125 sumptions, we present our main results with some simulations and discussions.
 126 Proofs are postponed till the next section.

127 Let the input space X be a separable metric space and ρ be a Borel proba-
 128 bility measure on $Z := X \times Y$.

For a predictor $f : X^2 \rightarrow \mathbb{R}$, we use a loss function $V : \mathbb{R}^2 \rightarrow \mathbb{R}_+$ and a
 reducing function $r : Y \times Y \rightarrow \mathbb{R}$ to give the local error $V(r(y, y'), f(x, x'))$ for
 $z = (x, y), z' = (x', y') \in Z$. The *generalization error* or risk $\mathcal{E} = \mathcal{E}^V$ associated
 with the loss function V is defined as

$$\mathcal{E}(f) = \int_Z \int_Z V(r(y, y'), f(x, x')) d\rho(z) d\rho(z').$$

129 We assume that there exists at least one minimizer f_ρ^V of the generalization
 130 error $\mathcal{E}(f)$, among all measurable functions $f : X^2 \rightarrow \mathbb{R}$. The goal of pairwise
 131 learning is to learn f_ρ^V from the sample set $S = \{z_t = (x_t, y_t)\}_{t=1}^T$ of size $T \in \mathbb{N}$.
 132 Throughout this paper, we assume that the samples are independently drawn
 133 according to ρ .

134 Our learning algorithm is a kernel method, where a RKHS is our hypothesis
 135 space. Let $K : X^2 \times X^2 \rightarrow \mathbb{R}$ be a Mercer Kernel, i.e., a continuous, symmetric
 136 and positive semi-definite kernel. The kernel K defines the RKHS $(\mathcal{H}_K, \|\cdot\|_K)$
 137 as the completion of the linear span of the set $\{K_{(x, x')}(\cdot) := K((x, x'), (\cdot, \cdot)) :$
 138 $(x, x') \in X^2\}$ with respect to an inner product $\langle \cdot, \cdot \rangle_K$ satisfying the reproducing
 139 property: i.e., $\langle K_{(x, x')}, g \rangle_K = g(x, x')$ for any $(x, x') \in X^2$ and $g \in \mathcal{H}_K$.

140 We assume in this paper that V is convex with respect to the second variable.
 141 That is, for any fixed $y \in \mathbb{R}$, the univariate function $V(y, \cdot)$ on \mathbb{R} is convex, hence
 142 its left-hand derivative $V'_-(y, f)$ exists at every $f \in \mathbb{R}$ and is non-decreasing.

143 The online pairwise learning algorithm considered in this paper is as follows.

144 **Algorithm 1.** *The online pairwise learning algorithm associated with the loss*

145 function V and the kernel K is defined by $f_1 = f_2 = 0$ and

$$f_{t+1} = f_t - \frac{\eta_t}{t-1} \sum_{j=1}^{t-1} V'_-(r(y_t, y_j), f_t(x_t, x_j)) K_{(x_t, x_j)}, \quad t = 2, \dots, T, \quad (2.1)$$

146 where $\{\eta_t > 0\}_t$ is a step size sequence.

147 The main purpose of this paper is to estimate the expected excess gener-
 148 alization error $\mathbb{E}[\mathcal{E}(f_T) - \mathcal{E}(f_\rho^V)]$. To this end, we shall make the following
 149 assumptions.

150 **Assumption 2.1.** *We assume*

$$|V|_0 := \sup_{y, y' \in Y} V(r(y, y'), 0) < \infty \quad (2.2)$$

151 and an increment condition for the left-hand derivative $V'_-(y, \cdot)$ that for some
 152 $q \geq 0$ and constant $c_q > 0$, there holds

$$|V'_-(r(y, y'), f)| \leq c_q(1 + |f|^q), \quad \forall f \in \mathbb{R}, y, y' \in Y. \quad (2.3)$$

153 We assume the kernel to be bounded with

$$\kappa = \max \left(\sup_{x, x' \in X} \sqrt{K((x, x'), (x, x'))}, 1 \right) < \infty. \quad (2.4)$$

154 Assumption (2.2) automatically holds for loss functions widely used for clas-
 155 sification, where V takes the form $V(y, f) = \phi(-yf)$ with $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ being
 156 a convex function, including the hinge loss ϕ_h , the exponential loss $\phi(v) =$
 157 $\exp(-v)$ and the logistic loss $\phi(v) = \log(1 + \exp(-v))$. Assumption (2.2) is
 158 equivalent to the boundedness assumption on the output space Y for $r(y, y') =$
 159 $y - y'$ and loss functions of the form $V(y, f) = \phi(y - f)$ for regression with
 160 $\lim_{|y| \rightarrow \infty} \phi(y) = \infty$, including the p -norm absolute distance loss $\phi(y) = |y|^p$
 161 with $p \geq 1$. Note that (2.2) may also hold for the case that Y is not bounded,
 162 e.g., the ranking problems with $r(y, y') = \text{sign}(y - y')$. The increment condition
 163 on loss functions (2.3) and the boundness assumption on the kernel are quite
 164 common in learning theory. For specific loss functions, one can easily compute
 165 the constants q and c_q in (2.3). For example, if the loss function is the hinge

166 loss $V(y, f) = \phi_h(yf)$, we know [25] that (2.3) is satisfied with $q = 0$ and
 167 $c_q = \sup_{y, y' \in Y} |r(y, y')|$, and in this case $|V|_0 = 1$.

168 We also need a notion of approximation error to state our main results.

169 **Definition 2.2.** *The approximation error associated with the tripe (ρ, V, K) is*
 170 *defined by*

$$\mathcal{D}(\lambda) = \inf_{f \in \mathcal{H}_K} \{ \mathcal{E}(f) - \mathcal{E}(f_\rho^V) + \lambda \|f\|_K^2 \}, \quad \forall \lambda > 0. \quad (2.5)$$

171 Our main result of this paper is stated as follows.

172 **Theorem 2.3.** *Under Assumption 2.1, let $\{\eta_{t+1} = \eta_1 t^{-\theta}\}_{t \in \mathbb{N}}$ with $\frac{q}{q+1} \leq \theta < 1$*
 173 *and η_1 satisfying*

$$0 < \eta_1 \leq \min \left\{ \frac{\sqrt{1-\theta}}{2\sqrt{2}c_q\kappa(\kappa+1)^q}, \frac{1-\theta}{4|V|_0} \right\}. \quad (2.6)$$

Then the sequence $\{f_t\}_t$ generated by Algorithm 1 satisfies

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} \leq \tilde{C}_0 \mathcal{D}((T-1)^{\theta-1}) + \tilde{C}_1 \Lambda_{T-1},$$

174 *where Λ_{T-1} is the quantity defined by*

$$\Lambda_{T-1} = \begin{cases} (T-1)^{-(1-\theta)}, & \text{when } \theta > \frac{q+2}{q+3}, \\ (T-1)^{-\frac{q\theta+\theta-q}{2}} \log(eT), & \text{when } \theta \leq \frac{q+2}{q+3}, \end{cases} \quad (2.7)$$

175 *and \tilde{C}_0 and \tilde{C}_1 are constants independent of T (given explicitly in the proof).*

176 To state explicit convergence rates, we need the following assumption for the
 177 decay of the approximation error.

178 **Assumption 2.4.** *Assume that for some $\beta \in (0, 1]$ and $c_\beta > 0$, the approxima-*
 179 *tion error satisfies*

$$\mathcal{D}(\lambda) \leq c_\beta \lambda^\beta, \quad \forall \lambda > 0. \quad (2.8)$$

180 The assumption (2.8) on the approximation error is independent of the sam-
 181 ples, and measures the approximation ability of the space \mathcal{H}_K to f_ρ^V with respect
 182 to (ρ, V) . It is standard in learning theory both in pairwise [32] and pointwise

183 learning [25, 29, 11]. Note that in the ideal case with $f_\rho^V \in \mathcal{H}_K$, condition (2.8)
 184 always holds with $\beta = 1$ and $c_\beta \leq \|f_\rho^V\|_K^2$.

185 We now have the following corollary, which follows directly from Theorem
 186 2.3.

187 **Corollary 2.5.** *Under the assumptions and notations of Theorem 2.3, and*
 188 *Assumption 2.4, we have*

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O(T^{(\theta-1)\beta} + \Lambda_T). \quad (2.9)$$

189 *In particular, we have*

$$(I) \text{ for } \theta = \frac{q+2}{q+3},$$

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O(T^{-\frac{\beta}{q+3}} \log T).$$

$$(II) \text{ for } \theta = \frac{q+2\beta}{q+1+2\beta},$$

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O(T^{-\frac{\beta}{q+1+2\beta}} \log T).$$

190 The above result gives bounds on the expected excess generalization error,
 191 where the general convergence rate in (2.9) depends on three parameters: q , β ,
 192 and θ . In general, it is easy to compute the increment parameter q for a given
 193 loss function, whereas the parameter β is unknown. Given only the growth
 194 parameter q , Part (I) of Corollary 2.5 suggests that the optimal convergence
 195 rate is achieved by setting $\theta = \frac{q+2}{q+3}$. If furthermore, the parameter β is provided,
 196 the optimal convergence rate is achieved by setting $\theta = \frac{q+2\beta}{q+1+2\beta}$.

197 Specifying the loss function in the above results, we have the following con-
 198 vergence rates with the hinge loss and the least squares loss.

199 **Corollary 2.6** (Hinge loss). *Let the loss function $V(y, a)$ be given with the hinge*
 200 *loss as $V(y, a) = \phi_h(ya)$. Assume (2.4), (2.8) and $M := \sup_{y, y' \in Y} |r(y, y')| <$*
 201 *∞ . Choose $\{\eta_{t+1} = \eta_1 t^{-\theta}\}_{t \in \mathbb{N}}$ with η_1 satisfying (2.6), where $q = 0$, $c_q = M$*
 202 *and $|V|_0 = 1$. Then for the sequence $\{f_t\}_t$ generated by Algorithm 1, we have*
 203 *the following convergence rates.*

(I) If $\theta = \frac{2}{3}$, then

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O \left(T^{-\frac{\beta}{3}} \log T \right).$$

204 Specially, if $\beta = 1$, i.e., $f_\rho^V \in \mathcal{H}_K$, then the upper bound is of order
 205 $O \left(T^{-\frac{1}{3}} \log T \right)$.

(II) If $\theta = \frac{2\beta}{2\beta+1}$, then

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O \left(T^{-\frac{\beta}{2\beta+1}} \log T \right).$$

206 **Corollary 2.7** (Least squares loss). *Let V be given by the least squares loss as*
 207 $V(y, a) = (y-a)^2$. *Assume (2.4), (2.8) and $M := 2 \max \left(\sup_{y, y' \in Y} |r(y, y')|, 1 \right) <$
 208 ∞ . *Choose $\{\eta_{t+1} = \eta_1 t^{-\theta}\}_{t \in \mathbb{N}}$ with η_1 satisfying (2.6), where $q = 1, c_q = M$ and*
 209 $|V|_0 = \sup_{y, y' \in Y} (r(y, y'))^2$. *Then for the sequence $\{f_t\}_t$ generated by Algorithm*
 210 *1, we have the following convergence rates.**

(I) If $\theta = \frac{3}{4}$, then

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O \left(T^{-\frac{\beta}{4}} \log T \right).$$

211 Specially, if $\beta = 1$, i.e., $f_\rho^V \in \mathcal{H}_K$, then the upper bound is of order
 212 $O \left(T^{-\frac{1}{4}} \log T \right)$.

(II) If $\theta = \frac{2\beta+1}{2\beta+2}$, then

$$\mathbb{E}_{z_1, \dots, z_T} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} = O \left(T^{-\frac{\beta}{2\beta+2}} \log T \right).$$

213 *Simulations..* We perform simulation experiments here to illustrate the derived
 214 convergence rates with polynomial decaying stepsizes. We consider the ranking
 215 problem with the loss function $V(y, a)$ given by the hinge loss as $V(y, a) =$
 216 $\phi_h(ya)$ and the reducing function $r(y, y') = \text{sign}(y - y')$. We consider the
 217 Boston housing dataset [13], which has 506 examples and 13 features, includ-
 218 ing *per capita crime rate by town, weighted distances to five Boston employ-*
 219 *ment centres and average number of rooms per dwelling.* We wish to predict
 220 the ordering based on values of houses and consider linear ranking rules with

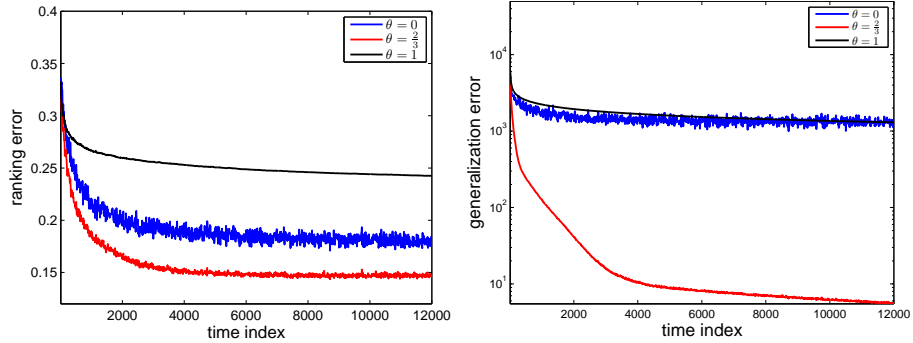


Figure 1: The behavior of Algorithm 1 on the Boston housing dataset. Left: ranking errors versus different stepsize sequences, right: generalization errors versus different stepsize sequences.

221 $K((x, x'), (u, u')) = (x - x')^\top (u - u')$ for $x, x', u, u' \in \mathbb{R}^{13}$. Here x^\top denotes
 222 the transpose of x . Let ρ be the uniform distribution on the 506 examples
 223 in the Boston housing dataset. We define the ranking error of a predictor
 224 $f : X \times X \rightarrow \mathbb{R}$ by $L(f) = \mathbb{E}[\text{sign}(y - y')f(x, x') < 0]$. We apply Algorithm 1
 225 with $\eta_t = (t - 1)^{-\theta}$ and $\theta \in \{0, 1, \frac{2}{3}\}$. We repeat the experiments 400 times and
 226 report the average ranking errors and average generalization errors. Figure 1
 227 illustrates the behavior of Algorithm 1 with three different stepsize sequences.
 228 It shows that the algorithm with polynomial decaying stepsize sequence with
 229 $\theta = \frac{2}{3}$ performs better than that with the constant stepsize sequence $\eta_t \equiv 1$
 230 and the sequence with $\theta = 1$. This is consistent with our theoretical results in
 231 Corollary 2.6.

232 *Discussions.* As mentioned before, online pairwise learning involves non-i.i.d.
 233 sample pairs. Thus, the analysis for pairwise learning is more challenging,
 234 in contrast with that for the online pointwise learning [6, 24, 31, 2, 22, 18].
 235 With the step size $\eta_t = \eta_1 t^{-\frac{\beta}{\beta+1}}$, the convergence rate $O(T^{-\frac{\beta}{\beta+1}} \log T)$ was
 236 established in [18] for the online pointwise learning, which is comparable to
 237 the convergence rate for batch learning in the pointwise setting. The con-
 238 vergence rate we derived in Corollary 2.5 for the online pairwise learning is
 239 of order $O(T^{-\frac{\beta}{2\beta+1+q}} \log T)$. This is due to an essential statistical difference

240 between these two families of learning algorithms: while the online pointwise
 241 learning uses unbiased estimators of the true gradients in the learning process,
 242 the randomized gradient $\frac{1}{t-1} \sum_{j=1}^{t-1} V'_-(r(y_t, y_j), f_t(x_t, x_j))K_{(x_t, x_j)}$ used in the
 243 online pairwise learning is a biased estimator of the true gradient $\int_Z \int_Z V'_-(y -$
 244 $y', f_t(x, x'))K_{(x, x')}d\rho(z)d\rho(z')$. We overcome this obstacle by applying the tool
 245 of Rademacher complexity to control the difference between partial generaliza-
 246 tion errors and generalization errors, resulting in, however, an additional term
 247 that dominates the upper bound in Proposition 3.6.

248 In what follows, we compare our work with existing results on online algo-
 249 rithms for pairwise learning. We first compare our work with [15, 27], where
 250 the online-to-batch conversion bounds for projected online pairwise learning
 251 algorithms in Euclidean spaces were provided. Assuming that $f_\rho^V \in \mathbb{R}^d$ is in the
 252 projected-bounded domain, upper bounds on the excess generalization error of
 253 order $O(T^{-\frac{1}{2}})$ were presented in [15] for the average iterates. In contrast, Algo-
 254 rithm 1 does not have any additional projection step and is implemented in the
 255 unconstrained setting on RKHSs including the Euclidean spaces. Besides, our
 256 bounds are stated in a more general setting for the last iterates, involving ap-
 257 proximation errors. It should be mentioned that convergence rates $O(T^{-\frac{1}{2}} \log T)$
 258 can be achieved by our analysis for the pairwise learning setting if an additional
 259 projection is performed at each iteration and $\beta = 1$. Finally, we compare our
 260 results with the existing work in [32, 33, 12]. Algorithm 1 with kernel methods
 261 was studied in [32] for the least squares loss, and in [33] for 1-activating loss V ,
 262 i.e., loss function which is differentiable and satisfies

$$|V'(y, f) - V'(y, g)| \leq L|f - g|, \quad \forall y \in \mathbb{R}, f, g \in \mathbb{R}, \quad (2.10)$$

263 for some $0 < L < \infty$. A convergence rate of order $O(T^{-\min\{\frac{\beta}{\beta+1}, \frac{1}{3}\}} \log T)$ is
 264 achieved for the algorithm with the least squares loss in [32]. However, the
 265 analysis in [32] is based on an integral operator approach and does not apply
 266 to general convex loss functions. Note that the results in [32] are in probability
 267 while our results are stated in expectation, and it would be interesting to further
 268 develop bounds in probability for the algorithm involving convex loss functions.

269 In comparison with the results in [33] where 1-activating loss functions are
 270 studied with an assumption on the existence of a minimizer of $\mathcal{E}(f)$ for $f \in$
 271 \mathcal{H}_K , our results hold for a broader class of loss functions and are better for
 272 1-activating loss functions in a more general setting. First, the hinge loss and
 273 the p -absolute value loss functions with $p \neq 2$ are not covered in [33]. Second,
 274 it is easy to see that an 1-activating loss function always satisfies the growth
 275 condition (2.3) with $q = 1$. Thus, by setting $\beta = 1$ and $\eta_t = \eta_1 t^{-\frac{\alpha+2}{\alpha+3}}$ in Corollary
 276 2.5, our optimal convergence rates are of order $O(T^{-\frac{1}{4}} \log T)$ for 1-activating loss
 277 functions, which are better than the bounds in [33] of order $O(T^{\epsilon-\frac{1}{6}})$ with an
 278 arbitrarily small $\epsilon > 0$. When the incremental exponent q satisfies $0 \leq q < 1$,
 279 the learning rates of order $O(T^{-\frac{\beta}{q+1+2\beta}} \log T)$ stated in Corollary 2.5 (II) are
 280 also better than those of order $O(T^{-\frac{\beta}{2\beta+2}} \sqrt{\log T})$ derived for online pairwise
 281 learning based on regularization schemes in RKHSs in the earlier work [12].

282 3. Proofs

283 In this section, we prove Theorem 2.3. To do so, it is necessary to prove
 284 some preliminary lemmas.

285 3.1. Bounding the learning sequence

For notational simplicity, we introduce the following two notations: the local
 empirical error of a function $f : X \times X \rightarrow \mathbb{R}$ at point z_t with respect to an ordered
 dataset $S = \{z_1, \dots, z_T\}$

$$\widehat{\mathcal{E}}_S^t(f) = \frac{1}{t-1} \sum_{j=1}^{t-1} V(r(y_t, y_j), f(x_t, x_j)),$$

and the partial generalization error with respect to an ordered dataset $S =$
 $\{z_1, \dots, z_T\}$

$$\widetilde{\mathcal{E}}_S^t(f) = \frac{1}{t-1} \sum_{j=1}^{t-1} \int_Z V(r(y, y_j), f(x, x_j)) d\rho(x, y).$$

286 We first introduce the following lemma whose proof essentially makes use of the
 287 convexity and increment property of loss functions.

288 **Lemma 3.1.** Under condition (2.3), for an arbitrary fixed $f \in \mathcal{H}_K$, and $t =$
 289 $2, \dots, T$,

$$\|f_{t+1} - f\|_K^2 \leq \|f_t - f\|_K^2 + \eta_t^2 G_t^2 + 2\eta_t(\widehat{\mathcal{E}}_S^t(f) - \widehat{\mathcal{E}}_S^t(f_t)), \quad (3.1)$$

290 where

$$G_t^2 = 4c_q^2 \kappa^2 (\kappa + 1)^{2q} \max \left\{ 1, \|f_t\|_K^{2q} \right\}. \quad (3.2)$$

Proof. Since f_{t+1} is given by (2.1), we have

$$\begin{aligned} \|f_{t+1} - f\|_K^2 &= \|f_t - f\|_K^2 + \eta_t^2 \left\| \frac{1}{t-1} \sum_{j=1}^{t-1} V'_-(r(y_t, y_j), f_t(x_t, x_j)) K_{(x_t, x_j)} \right\|_K^2 \\ &\quad + \frac{2\eta_t}{t-1} \sum_{j=1}^{t-1} V'_-(r(y_t, y_j), f_t(x_t, x_j)) \langle K_{(x_t, x_j)}, f - f_t \rangle_K. \end{aligned}$$

Observe that

$$\|K_{(x_t, x_j)}\|_K = \sqrt{K((x_t, x_j), (x_t, x_j))} \leq \kappa$$

and that

$$\|f\|_\infty \leq \kappa \|f\|_K, \quad \forall f \in \mathcal{H}_K.$$

291 These together with the increment condition (2.3) yield

$$\begin{aligned} \|V'_-(r(y_t, y_j), f_t(x_t, x_j)) K_{(x_t, x_j)}\|_K &\leq \kappa |V'_-(r(y_t, y_j), f_t(x_t, x_j))| \\ &\leq \kappa c_q (1 + |f_t(x_t, x_j)|^q) \leq \kappa c_q (1 + \kappa^q \|f_t\|_K^q). \end{aligned}$$

Therefore,

$$\|f_{t+1} - f\|_K^2 \leq \|f_t - f\|_K^2 + \eta_t^2 G_t^2 + \frac{2\eta_t}{t-1} \sum_{j=1}^{t-1} V'_-(r(y_t, y_j), f_t(x_t, x_j)) \langle K_{(x_t, x_j)}, f - f_t \rangle_K.$$

292 Using the reproducing property, we get

$$\|f_{t+1} - f\|_K^2 \leq \|f_t - f\|_K^2 + \eta_t^2 G_t^2 + \frac{2\eta_t}{t-1} \sum_{j=1}^{t-1} V'_-(r(y_t, y_j), f_t(x_t, x_j)) (f(x_t, x_j) - f_t(x_t, x_j)). \quad (3.3)$$

Since $V(r(y_t, y_j), \cdot)$ is a convex function, we have

$$V'_-(r(y_t, y_j), a)(b - a) \leq V(r(y_t, y_j), b) - V(r(y_t, y_j), a), \quad \forall a, b \in \mathbb{R}.$$

293 Using this relation in (3.3), we get our desired result. \square

294 Using the above lemma, we can bound the learning sequence as follows.
 295 The proof is motivated by the recent work in [16] and [17], using an induction
 296 argument.

297 **Lemma 3.2.** *Assume condition (2.3). Let $\frac{q}{q+1} \leq \theta < 1$ and $\eta_{t+1} = \eta_1 t^{-\theta}$ for
 298 $t \in \mathbb{N}$ with η_1 satisfying (2.6). Then for $t = 1, \dots, T$,*

$$\|f_{t+1}\|_K \leq (t-1)^{\frac{1-\theta}{2}}. \quad (3.4)$$

299 *Proof.* We prove our statement by induction.

Taking $f = 0$ in Lemma 3.1, we know that

$$\|f_{t+1}\|_K^2 \leq \|f_t\|_K^2 + \eta_t^2 G_t^2 + 2\eta_t(\widehat{\mathcal{E}}_S^t(0) - \widehat{\mathcal{E}}_S^t(f_t)) \leq \|f_t\|_K^2 + \eta_t^2 G_t^2 + 2\eta_t|V|_0.$$

300 This verifies (3.4) for the case $t = 2$ since $f_1 = f_2 = 0$ and $4\eta_1^2 c_q^2 \kappa^2 (\kappa + 1)^{2q} +$
 301 $2\eta_1|V|_0 \leq 1$.

302 Assume $\|f_t\|_K \leq (t-2)^{\frac{1-\theta}{2}}$ with $t \geq 3$. Then

$$G_t^2 \leq 4c_q^2 \kappa^2 (\kappa + 1)^{2q} (t-2)^{(1-\theta)q}. \quad (3.5)$$

303 Hence

$$\begin{aligned} \|f_{t+1}\|_K^2 &\leq (t-2)^{1-\theta} + 4\eta_1^2 (t-1)^{-2\theta} c_q^2 \kappa^2 (\kappa + 1)^{2q} (t-1)^{(1-\theta)q} + 2\eta_1 (t-1)^{-\theta} |V|_0 \\ &\leq (t-1)^{1-\theta} \left\{ \left(1 - \frac{1}{t-1}\right)^{1-\theta} + \frac{4\eta_1^2 c_q^2 \kappa^2 (\kappa + 1)^{2q}}{(t-1)^{(q+1)\theta+1-q}} + \frac{2\eta_1 |V|_0}{t-1} \right\}. \end{aligned}$$

Since $\left(1 - \frac{1}{t-1}\right)^{1-\theta} \leq 1 - \frac{1-\theta}{t-1}$ and the condition $\theta \geq \frac{q}{q+1}$ implies $(q+1)\theta+1-q \geq 1$, we have

$$\|f_{t+1}\|_K^2 \leq (t-1)^{1-\theta} \left\{ 1 - \frac{1-\theta}{t-1} + \frac{4\eta_1^2 c_q^2 \kappa^2 (\kappa + 1)^{2q}}{t-1} + \frac{2\eta_1 |V|_0}{t-1} \right\}.$$

304 Finally we use the restriction (2.6) for η_1 and find $\|f_{t+1}\|_K^2 \leq (t-1)^{1-\theta}$. This
 305 completes the induction procedure and proves our conclusion. \square

306 With the above two lemmas, and noticing that f_t is independent of z_t , we
 307 can easily prove the following result.

308 **Proposition 3.3.** *Assume condition (2.3). Let $\frac{q}{q+1} \leq \theta < 1$ and $\eta_{t+1} = \eta_1 t^{-\theta}$
309 for all $t \in \mathbb{N}$ with η_1 satisfying (2.6). Assume that $t \in \{2, \dots, T\}$ and that
310 $f \in \mathcal{H}_K$ is independent of z_t (but may depend on z_1, \dots, z_{t-1}). Then we have*

$$\begin{aligned} \mathbb{E}_{z_t} \|f_{t+1} - f\|_K^2 &\leq \|f_t - f\|_K^2 \\ &+ 4\eta_1^2 c_q^2 \kappa^2 (\kappa + 1)^{2q} (t-1)^{(1-\theta)q-2\theta} + 2\eta_t \left[\tilde{\mathcal{E}}_S^t(f) - \tilde{\mathcal{E}}_S^t(f_t) \right]. \end{aligned} \quad (3.6)$$

Proof. Taking expectations on both sides of (3.1) with respect to z_t , and noting that f_t is independent of z_t , we get

$$\mathbb{E}_{z_t} \|f_{t+1} - f\|_K^2 \leq \|f_t - f\|_K^2 + \eta_t^2 G_t^2 + 2\eta_t \left[\tilde{\mathcal{E}}_S^t(f) - \tilde{\mathcal{E}}_S^t(f_t) \right].$$

311 Lemma 3.2 shows that $\|f_t\|_K \leq (t-1)^{\frac{1-\theta}{2}}$, which implies (3.5). Applying (3.5)
312 and using $\eta_t = \eta_1 (t-1)^{-\theta}$ in the above inequality yield the desired bound. \square

313 Proposition 3.3 gives an iterated inequality related to the partial generaliza-
314 tion error $\tilde{\mathcal{E}}_S^t(f_t)$. Note that our goal is to derive upper bounds on the excess
315 generalization error. It is thus necessary to develop relationships between the
316 partial generalization error and generalization error, which will be considered in
317 the following subsection.

318 3.2. From partial generalization error to generalization error

319 For $R > 0$, denote B_R the ball of radius R in \mathcal{H}_K : $B_R = \{f \in \mathcal{H}_K : \|f\|_K \leq$
320 $R\}$. The following lemma gives a uniform upper bound on the differences be-
321 tween the partial generalization error and generalization error over any ball B_R
322 with $R \geq 1$. Its proof uses a standard symmetry technique and some properties
323 related to Rademacher complexity.

Lemma 3.4. *For $R \geq 1$, and all $1 \leq t \leq T$*

$$\mathbb{E}_{z_1, \dots, z_{t-1}} \left[\sup_{f \in B_R} \{\mathcal{E}(f) - \tilde{\mathcal{E}}_S^t(f)\} \right] \leq \frac{2c_q R \kappa (1 + \kappa^q R^q)}{\sqrt{t-1}}.$$

324 *The above inequality is also true if we replace $\{\mathcal{E}(f) - \tilde{\mathcal{E}}_S^t(f)\}$ by $\{\tilde{\mathcal{E}}_S^t(f) - \mathcal{E}(f)\}$.*

Proof. For notational simplicity, we denote

$$\mathcal{L}(f, z_j) = \int_Z V(r(y, y_j), f(x, x_j)) d\rho(z).$$

Then

$$\tilde{\mathcal{E}}_S^t(f) = \frac{1}{t-1} \sum_{j=1}^{t-1} \mathcal{L}(f, z_j)$$

and

$$\mathcal{E}(f) = \int_Z \mathcal{L}(f, z') d\rho(z').$$

325 Let $S' = \{z'_1, \dots, z'_T\}$ be another independent sample set. We first note that

$$\begin{aligned} & \mathbb{E}_S[\sup_{f \in B_R} \{\mathcal{E}(f) - \tilde{\mathcal{E}}_S^t(f)\}] \\ &= \mathbb{E}_S[\sup_{f \in B_R} \{\mathbb{E}_{S'}[\tilde{\mathcal{E}}_{S'}^t(f)] - \tilde{\mathcal{E}}_S^t(f)\}] \\ &\leq \mathbb{E}_{S, S'}[\sup_{f \in B_R} \{\tilde{\mathcal{E}}_{S'}^t(f) - \tilde{\mathcal{E}}_S^t(f)\}]. \end{aligned}$$

326 Here, we abuse the notation \mathbb{E}_S for $\mathbb{E}_{z_1, \dots, z_{t-1}}$. Let $\sigma_1, \sigma_2, \dots, \sigma_T$ be independent
 327 random variables drawn from the Rademacher distribution i.e. $\Pr(\sigma_i = +1) =$
 328 $\Pr(\sigma_i = -1) = 1/2$ for $i = 1, 2, \dots, T$. Using a standard symmetry technique,
 329 for example in [3],

$$\begin{aligned} & \mathbb{E}_{S, S'}[\sup_{f \in B_R} \{\tilde{\mathcal{E}}_{S'}^t(f) - \tilde{\mathcal{E}}_S^t(f)\}] \\ &\leq \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in B_R} \left\{ \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j (\mathcal{L}(f, z'_j) - \mathcal{L}(f, z_j)) \right\} \right]. \end{aligned}$$

330 Thus,

$$\begin{aligned}
& \mathbb{E}_S[\sup_{f \in B_R} \{\mathcal{E}(f) - \tilde{\mathcal{E}}_S^t(f)\}] \\
& \leq \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in B_R} \left\{ \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j (\mathcal{L}(f, z'_j) - \mathcal{L}(f, z_j)) \right\} \right] \\
& \leq 2\mathbb{E}_{S, \sigma} \left[\sup_{f \in B_R} \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j \mathcal{L}(f, z_j) \right] \\
& = 2\mathbb{E}_{S, \sigma} \left[\sup_{f \in B_R} \mathbb{E}_z \left[\frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j V(r(y, y_j), f(x, x_j)) \right] \right] \\
& \leq 2\mathbb{E}_{z, S, \sigma} \left[\sup_{f \in B_R} \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j V(r(y, y_j), f(x, x_j)) \right].
\end{aligned}$$

For any $z \in Z$, the term $\mathbb{E}_{S, \sigma} \left[\sup_{f \in B_R} \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j V(r(y, y_j), f(x, x_j)) \right]$ is the Rademacher complexity [4] of the function class B_R with respect to ρ for sample size $t-1$. Using (2.3) and that $\|f\|_\infty \leq \kappa \|f\|_K$, it is easy to see that for any $f, f' \in B_R$,

$$|V(r(y, y_j), f(x, x_j)) - V(r(y, y_j), f'(x, x_j))| \leq c_q(1 + R^q \kappa^q) |f(x, x_j) - f'(x, x_j)|.$$

331 Applying Talagrand's contraction lemma (see e.g., [19, Theorem 7]), we have

$$\begin{aligned}
& \mathbb{E}_{S, \sigma} \left[\sup_{f \in B_R} \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j V(r(y, y_j), f(x, x_j)) \right] \\
& \leq c_q(1 + \kappa^q R^q) \mathbb{E}_{S, \sigma} \left[\sup_{f \in B_R} \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j f(x, x_j) \right]
\end{aligned}$$

332 and therefore,

$$\begin{aligned}
& \mathbb{E}_S[\sup_{f \in B_R} \mathbb{E}\{\mathcal{E}(f) - \tilde{\mathcal{E}}^t(f)\}] \\
& \leq 2c_q(1 + \kappa^q R^q) \mathbb{E}_{z, S, \sigma} \left[\sup_{f \in B_R} \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j f(x, x_j) \right] \\
& = 2c_q(1 + \kappa^q R^q) \mathbb{E}_{z, S, \sigma} \left[\sup_{f \in B_R} \left\langle f, \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j K_{(x, x_j)} \right\rangle_K \right].
\end{aligned}$$

333 Applying the Schwarz inequality,

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in B_R} \mathbb{E} \{ \mathcal{E}(f) - \tilde{\mathcal{E}}^t(f) \} \right] \\ & \leq 2c_q (1 + \kappa^q R^q) \mathbb{E}_{z, S, \sigma} \left[\sup_{f \in B_R} \|f\|_K \left\| \frac{1}{t-1} \sum_{j=1}^{t-1} \sigma_j K(x, x_j) \right\|_K \right]. \end{aligned}$$

334 Applying $\mathbb{E}[\|g\|_K] \leq (\mathbb{E}[\|g\|_K^2])^{\frac{1}{2}}$, and noting that $\sigma_1, \sigma_2, \dots, \sigma_T$ are independent
335 random variables with mean zeros,

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in B_R} \mathbb{E} \{ \mathcal{E}(f) - \tilde{\mathcal{E}}^t(f) \} \right] \\ & \leq \frac{2c_q (1 + \kappa^q R^q) R}{t-1} \left[\mathbb{E}_{z, S, \sigma} \left\| \sum_{j=1}^{t-1} \sigma_j K(x, x_j) \right\|_K^2 \right]^{\frac{1}{2}} \\ & = \frac{2c_q (1 + \kappa^q R^q) R}{t-1} \left[\sum_{j=1}^{t-1} \mathbb{E}_{x, x_j} \|K(x, x_j)\|_K^2 \right]^{\frac{1}{2}} \\ & \leq \frac{2c_q (1 + \kappa^q R^q) R \kappa}{\sqrt{t-1}}, \end{aligned}$$

336 where for the last inequality we use the boundness assumption on the kernel.
337 Thus we get the desired result. The proof is complete. \square

338 Combining the above lemma with Lemma 3.2, we get the following corollary.

Corollary 3.5. *Under the assumptions of Lemma 3.2, we have for any $t = 3, \dots, T$,*

$$|\mathbb{E}_{z_1, \dots, z_{t-1}} [\mathcal{E}(f_t) - \tilde{\mathcal{E}}_S^t(f_t)]| \leq 2c_q \kappa (1 + \kappa^q) (t-1)^{\frac{(1-\theta)(q+1)-1}{2}}.$$

339 *3.3. A useful proposition*

340 The following proposition will be used several times in our proof. Its proof
341 follows directly from Proposition 3.3 and Corollary 3.5.

342 **Proposition 3.6.** *Under assumptions of Proposition 3.3, for any $f \in \mathcal{H}_K$*
343 *which is independent of z_1, \dots, z_t , or $f = f_k$ ($3 \leq k \leq t$), we have*

$$\begin{aligned} & 2\eta_t \mathbb{E}_{z_1, \dots, z_{t-1}} [\mathcal{E}(f_t) - \mathcal{E}(f)] \\ & \leq \mathbb{E}_{z_1, \dots, z_t} \{ \|f_t - f\|_K^2 - \|f_{t+1} - f\|_K^2 \} + C_{q, \kappa, \eta_1} (t-1)^{-q^*}. \end{aligned} \tag{3.7}$$

344 Here

$$q^* = \frac{3\theta - (1 - \theta)q}{2}. \quad (3.8)$$

345 and C_{q,κ,η_1} is a constant depending only on q, κ and η_1 , given explicitly by (3.10)
 346 in the proof.

347 *Proof.* Note that for $3 \leq k \leq T$, f_k depends only on z_1, \dots, z_{k-1} . By Proposi-
 348 tion 3.3, we have

$$\begin{aligned} \mathbb{E}_{z_1, \dots, z_t} \|f_{t+1} - f\|_K^2 &\leq \mathbb{E}_{z_1, \dots, z_t} \|f_t - f\|_K^2 \\ &+ 4\eta_1^2 c_q^2 \kappa^2 (\kappa + 1)^{2q} (t-1)^{(1-\theta)q-2\theta} + 2\eta_t \mathbb{E}_{z_1, \dots, z_{t-1}} \left[\tilde{\mathcal{E}}_S^t(f) - \tilde{\mathcal{E}}_S^t(f_t) \right]. \end{aligned}$$

349 Rewrite $\mathbb{E}_{z_1, \dots, z_{t-1}} \left[\tilde{\mathcal{E}}_S^t(f) - \tilde{\mathcal{E}}_S^t(f_t) \right]$ as

$$\mathbb{E}_{z_1, \dots, z_{t-1}} [\mathcal{E}(f) - \mathcal{E}(f_t)] + \mathbb{E}_{z_1, \dots, z_{t-1}} \left[(\tilde{\mathcal{E}}_S^t(f) - \mathcal{E}(f)) + (\mathcal{E}(f_t) - \tilde{\mathcal{E}}_S^t(f_t)) \right]. \quad (3.9)$$

If $f = f_k$ with $1 \leq k \leq t$, by applying Corollary 3.5 to bound the last term of
 (3.9), and noting that $\theta \geq \frac{q}{q+1}$ implies

$$\frac{(1-\theta)(q+1) - 1}{2} - \theta = \frac{(1-\theta)q - 3\theta}{2} \geq (1-\theta)q - 2\theta,$$

350 we get (3.7) with

$$C_{q,\kappa,\eta_1} = 4\eta_1^2 c_q^2 \kappa^2 (\kappa + 1)^{2q} + 8\eta_1 c_q \kappa (1 + \kappa^q). \quad (3.10)$$

If f is independent of z_1, \dots, z_t , the last term of (3.9) is exactly

$$\mathbb{E}_{z_1, \dots, z_{t-1}} \left[\mathcal{E}(f_t) - \tilde{\mathcal{E}}_S^t(f_t) \right].$$

351 Using Corollary 3.5 to bound this term again, we get (3.7). From the above
 352 analysis, one can conclude the proof. \square

353 3.4. Estimating excess generalization error

354 We now give the following general result, with which we can prove our main
 355 result, Theorem 2.3. For notational simplicity, we denote the excess generaliza-
 356 tion error of $f_* \in \mathcal{H}_K$ with respect to (ρ, V) by $\mathcal{A}(f_*)$:

$$\mathcal{A}(f_*) = \mathcal{E}(f_*) - \mathcal{E}(f_\rho^V). \quad (3.11)$$

357 **Theorem 3.7.** Assume (2.3) with $q \geq 0$. Let $\eta_{t+1} = \eta_1 t^{-\theta}$ with $\frac{q}{q+1} \leq \theta < 1$
 358 and η_1 satisfying (2.6). Then for every fixed $f_* \in \mathcal{H}_K$,

$$\mathbb{E}_{z_1, \dots, z_{T-1}} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} \leq \frac{\mathcal{A}(f_*)}{1-\theta} + \frac{\|f_*\|_K^2}{2\eta_1} (T-1)^{\theta-1} + \tilde{C}_1 \Lambda_{T-1}, \quad (3.12)$$

359 where Λ_{T-1} is given by (2.7) and \tilde{C}_1 is a positive constant depending on q, κ, θ
 360 (independent of T and f_* , and given explicitly in the proof).

361 The proof of this theorem is postponed to the next subsection. A novel
 362 error decomposition plays an important role in the proof. Note that the de-
 363 composition of ρ into the margin probability measure on X and the conditional
 364 probability measures allows the case with noise.

365 Now we are in a position to prove Theorem 2.3.

Proof of Theorem 2.3. By Theorem 3.7, we have

$$\mathbb{E}_{z_1, \dots, z_{T-1}} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} \leq \tilde{C}_0 (\mathcal{E}(f_*) - \mathcal{E}(f_\rho^V)) + (T-1)^{\theta-1} \|f_*\|_K^2 + \tilde{C}_1 \Lambda_{T-1},$$

where

$$\tilde{C}_0 = \frac{1}{1-\theta} + \frac{1}{2\eta_1}.$$

366 Since the constants \tilde{C}_0 and \tilde{C}_1 are independent of $f_* \in \mathcal{H}_K$, we can take infimum
 367 over $f_* \in \mathcal{H}_K$ on both sides, and conclude the desired result. \square

368 3.5. Proof of Theorem 3.7

369 Before proving Theorem 3.7, we present two lemmas, whose proofs follow
 370 from Proposition 3.6 and some elementary inequalities. In the rest of this sub-
 371 section, we denote $\mathbb{E}_{z_1, \dots, z_T}$ by \mathbb{E} for simplicity.

372 **Lemma 3.8** (Weighted average). *Under the assumption of Theorem 3.7, for*
 373 *any $T \geq 2$,*

$$\begin{aligned} \frac{1}{T-1} \sum_{t=2}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_\rho^V) \} &\leq \frac{\|f_*\|_K^2}{T-1} + \frac{2\eta_1 \mathcal{A}(f_*)}{1-\theta} (T-1)^{-\theta} \\ &+ \begin{cases} \frac{q^* C_{q, \kappa, \eta_1}}{q^* - 1} (T-1)^{-1}, & \text{when } \theta > \frac{q+2}{q+3}, \\ C_{q, \kappa, \eta_1} (T-1)^{-1} \log(eT), & \text{when } \theta = \frac{q+2}{q+3}, \\ \frac{C_{q, \kappa, \eta_1}}{1-q^*} (T-1)^{-q^*}, & \text{when } \theta < \frac{q+2}{q+3}. \end{cases} \end{aligned}$$

374 Here q^* and C_{q,κ,η_1} are given by (3.8) and (3.10), respectively.

375 *Proof.* Note that by Proposition 3.6, we have (3.7). Choosing $f = f_*$ in (3.7)
 376 and adding both sides with $2\eta_t\mathcal{A}(f_*)$, we get

$$\begin{aligned} & 2\eta_t\mathbb{E}[\mathcal{E}(f_t) - \mathcal{E}(f_\rho^V)] \\ & \leq \mathbb{E}\{\|f_t - f_*\|_K^2 - \|f_{t+1} - f_*\|_K^2\} + C_{q,\kappa,\eta_1}(t-1)^{-q^*} + 2\eta_t\mathcal{A}(f_*), \end{aligned}$$

377 Taking summations over $t = 2, \dots, T$, with $f_2 = 0$, and $\eta_t = \eta_1(t-1)^{-\theta}$,

$$\sum_{t=2}^T 2\eta_t\mathbb{E}\{\mathcal{E}(f_t) - \mathcal{E}(f_\rho^V)\} \leq \|f_*\|_K^2 + C_{q,\kappa,\eta_1} \sum_{t=1}^{T-1} t^{-q^*} + 2\eta_1\mathcal{A}(f_*) \sum_{t=1}^{T-1} t^{-\theta}.$$

Note that q^* is given by (3.8), and that from the restriction $\theta \in [\frac{q}{q+1}, 1)$, q^* satisfies $0 < q^* < 2$ and

$$q^* \begin{cases} > 1 & \text{when } \theta > \frac{q+2}{q+3}, \\ = 1 & \text{when } \theta = \frac{q+2}{q+3}, \\ < 1 & \text{when } \theta < \frac{q+2}{q+3}. \end{cases}$$

378 Applying

$$\sum_{t=1}^{T-1} t^{-\theta'} \leq 1 + \int_1^{T-1} u^{-\theta'} du \leq \begin{cases} \frac{(T-1)^{1-\theta'}}{1-\theta'}, & \text{when } \theta' < 1, \\ \log(eT), & \text{when } \theta' = 1, \\ \frac{\theta'}{\theta'-1}, & \text{when } \theta' > 1, \end{cases} \quad (3.13)$$

379 to bound $\sum_{t=1}^{T-1} t^{-q^*}$ and $\sum_{t=1}^{T-1} t^{-\theta}$, we get the desired result. The proof is
 380 complete. \square

381 **Lemma 3.9** (Moving weighted average). *Under the assumption of Theorem*
 382 *3.7, for any $T \geq 2$,*

$$\begin{aligned} & \sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T 2\eta_t\mathbb{E}\{\mathcal{E}(f_t) - \mathcal{E}(f_{T-k})\} \\ & \leq \begin{cases} 2C_{q,\kappa,\eta_1} \left(2^{q^*} + \frac{q^*}{q^*-1}\right) (T-1)^{-1}, & \text{when } \theta > \frac{q+2}{q+3}, \\ 4C_{q,\kappa,\eta_1} (\log T)(T-1)^{-1}, & \text{when } \theta = \frac{q+2}{q+3}, \\ 2C_{q,\kappa,\eta_1} \left(2^{q^*} + \frac{1}{1-q^*}\right) (\log T)(T-1)^{-q^*}, & \text{when } \theta < \frac{q+2}{q+3}. \end{cases} \end{aligned}$$

383 Here q^* and C_{q,κ,η_1} are given by (3.8) and (3.10), respectively.

384 *Proof.* Let $k \in \{2, \dots, T-1\}$. Note that f_{T-k} depends only on z_1, \dots, z_{T-k-1} .

385 By Proposition 3.6, we have for $t \geq T-k$,

$$2\eta_t \mathbb{E} [\mathcal{E}(f_t) - \mathcal{E}(f)] \leq \mathbb{E} \{ \|f_t - f\|_K^2 - \|f_{t+1} - f\|_K^2 \} + C_{q,\kappa,\eta_1} (t-1)^{-q^*}.$$

386 Taking summation over $t = T-k, \dots, T$ yields

$$\begin{aligned} \sum_{t=T-k+1}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_{T-k}) \} &= \sum_{t=T-k}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_{T-k}) \} \\ &\leq C_{q,\kappa,\eta_1} \sum_{t=T-k}^T (t-1)^{-q^*} = C_{q,\kappa,\eta_1} \sum_{t=T-1-k}^{T-1} t^{-q^*}. \end{aligned}$$

It thus follows that

$$\sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_{T-k}) \} \leq C_{q,\kappa,\eta_1} \sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{t=T-1-k}^{T-1} t^{-q^*}.$$

387 By applying the following elementary inequality from [16] (which will be proved

388 in the appendix for completeness)

$$\sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k}^T t^{-q^*} \leq \begin{cases} 2 \left(2^{q^*} + \frac{q^*}{q^*-1} \right) T^{-1}, & \text{when } q^* \in (1, 2), \\ 4(\log T) T^{-1}, & \text{when } q^* = 1, \\ 2 \left(2^{q^*} + \frac{1}{1-q^*} \right) (\log T) T^{-q^*}, & \text{when } q^* \in (0, 1), \end{cases} \quad (3.14)$$

389 the desired estimate is verified. The proof is complete. \square

390 With the above two lemmas, now we are ready to prove Theorem 3.7.

391 *Proof of Theorem 3.7.* The basic idea is to bound the weighted excess gener-

392 alization error $2\eta_T \mathbb{E}_{z_1, \dots, z_{T-1}} [\mathcal{E}(f_T) - \mathbb{E}(f_\rho^V)]$ in terms of the weighted average

393 and the moving weighted average. To do so, we need the following fact from

394 [22, 18] which asserts that for any sequence $\{u_j\}_{j \in \mathbb{N}}$ in \mathbb{R} , there holds

$$u_T = \frac{1}{T-1} \sum_{j=2}^T u_j + \sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{j=T-k+1}^T (u_j - u_{T-k}). \quad (3.15)$$

395 In fact, for $k \in \{1, \dots, T-2\}$, we have

$$\begin{aligned}
& \frac{1}{k} \sum_{j=T-k+1}^T u_j - \frac{1}{k+1} \sum_{j=T-k}^T u_j \\
&= \frac{1}{k(k+1)} \left\{ (k+1) \sum_{j=T-k+1}^T u_j - k \sum_{j=T-k}^T u_j \right\} \\
&= \frac{1}{k(k+1)} \sum_{j=T-k+1}^T (u_j - u_{T-k}).
\end{aligned}$$

396 Summing over $k = 2, \dots, T-1$, and rearranging terms, we get (3.15). Now, for

397 any $k = 1, \dots, T-2$, we choose $u_t = 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_\rho^V) \}$ in (3.15) to get

$$\begin{aligned}
2\eta_T \mathbb{E} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} &= \frac{1}{T-1} \sum_{j=2}^T 2\eta_j \mathbb{E} \{ \mathcal{E}(f_j) - \mathcal{E}(f_\rho^V) \} \\
&+ \sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{j=T-k+1}^T \left(2\eta_j \mathbb{E} \{ \mathcal{E}(f_j) - \mathcal{E}(f_\rho^V) \} - 2\eta_{T-k} \mathbb{E} \{ \mathcal{E}(f_{T-k}) - \mathcal{E}(f_\rho^V) \} \right),
\end{aligned}$$

398 which can be rewritten as

$$\begin{aligned}
2\eta_T \mathbb{E} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} &= \frac{1}{T-1} \sum_{t=2}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_\rho^V) \} \\
&+ \sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_{T-k}) \} \\
&+ \sum_{k=1}^{T-2} \frac{1}{k+1} \left[\frac{1}{k} \sum_{t=T-k+1}^T 2\eta_t - 2\eta_{T-k} \right] \mathbb{E} \{ \mathcal{E}(f_{T-k}) - \mathcal{E}(f_\rho^V) \}. \quad (3.16)
\end{aligned}$$

399 Since, $\mathcal{E}(f_{T-k}) - \mathcal{E}(f_\rho^V) \geq 0$ and that $\{\eta_t\}_{t \in \mathbb{N}}$ is a non-increasing sequence, we

400 know that the last term of the above inequality is at most zero. Therefore, we

401 get

$$\begin{aligned}
2\eta_T \mathbb{E} \{ \mathcal{E}(f_T) - \mathcal{E}(f_\rho^V) \} &\leq \frac{1}{T-1} \sum_{t=2}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_\rho^V) \} \\
&+ \sum_{k=1}^{T-2} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T 2\eta_t \mathbb{E} \{ \mathcal{E}(f_t) - \mathcal{E}(f_{T-k}) \}. \quad (3.17)
\end{aligned}$$

Applying lemmas 3.8 and 3.9 to bound the last two terms, we get the desired

bound (3.12) with \tilde{C}_1 given explicitly by

$$\tilde{C}_1 = \begin{cases} \frac{C_{q,\kappa,\eta_1}(3q^*+2^{q^*+1}(q^*-1))}{2\eta_1(q^*-1)}, & \text{when } \theta > \frac{q+2}{q+3}, \\ \frac{3C_{q,\kappa,\eta_1}}{\eta_1}, & \text{when } \theta = \frac{q+2}{q+3}, \\ \frac{C_{q,\kappa,\eta_1}\left(2^{q^*+1}+\frac{3}{1-q^*}\right)}{2\eta_1}, & \text{when } \theta < \frac{q+2}{q+3}. \end{cases}$$

402 The proof of Theorem 3.7 is complete. □

403 4. Conclusion

404 This paper presents learning rates of the last iterate for online pairwise learn-
 405 ing algorithms involving general convex loss functions which are better than the
 406 existing results under certain circumstances. Our idea is to use an error decom-
 407 position from [16, 23] to decompose the weighted excess generalization error into
 408 weighted average errors and moving weighted average errors. We apply some
 409 tools from Rademacher complexity to overcome the difficulty with the bias of the
 410 randomized gradients as estimators of the true gradients in the online pairwise
 411 learning setting. It is interesting to discuss here the connection between classifi-
 412 cation/regression tasks and pairwise learning problems. For the specific pairwise
 413 learning problem with $V(y, f) = (y - f)^2$ and $r(y, y') = y - y'$, it was proved
 414 in [32, 10] that the optimal predictor is $f_\rho^V(x, x') = \int_X y d\rho(y|x) - \int_X y d\rho(y|x')$,
 415 where $\rho(y|x)$ is the conditional measure at x . This shows that the pairwise learn-
 416 ing based on the least squares loss is essentially a pointwise learning problem
 417 since $\tilde{f}_\rho(x) := \int_X y d\rho(y|x)$ is the regression function minimizing $\int_Z (y - f(x))^2 d\rho$.
 418 Characterizing f_ρ^V and the approximation error assumption (2.8) for a general
 419 pairwise learning loss function in terms of function space properties, such as for
 420 metric and similarity learning, is a challenging problem for further study.

421 References

- 422 [1] S. Agarwal, P. Niyogi, Generalization bounds for ranking algorithms via
 423 algorithmic stability, *The Journal of Machine Learning Research* 10 (2009)
 424 441–474.

- 425 [2] F. Bach, E. Moulines, Non-strongly-convex smooth stochastic approxima-
426 tion with convergence rate $O(1/n)$, in: Advances in Neural Information
427 Processing Systems, 773–781, 2013.
- 428 [3] P. L. Bartlett, O. Bousquet, S. Mendelson, Local rademacher complexities,
429 Annals of Statistics (2005) 1497–1537.
- 430 [4] P. L. Bartlett, S. Mendelson, Rademacher and Gaussian complexities: Risk
431 bounds and structural results, The Journal of Machine Learning Research
432 3 (2003) 463–482.
- 433 [5] Q. Cao, Z.-C. Guo, Y. Ying, Generalization bounds for metric and similar-
434 ity learning, Machine Learning 102 (1) (2016) 115–132.
- 435 [6] N. Cesa-Bianchi, A. Conconi, C. Gentile, On the generalization ability
436 of on-line learning algorithms, IEEE Transactions on Information Theory
437 50 (9) (2004) 2050–2057, ISSN 0018-9448.
- 438 [7] A. Christmann, D.-X. Zhou, On the robustness of regularized pairwise
439 learning methods based on kernels, Journal of Complexity 37 (2016) 1–
440 33.
- 441 [8] S. Cléménçon, G. Lugosi, N. Vayatis, Ranking and empirical minimization
442 of U-statistics, The Annals of Statistics (2008) 844–874.
- 443 [9] F. Cucker, D.-X. Zhou, Learning Theory: an Approximation Theory View-
444 point, vol. 24, Cambridge University Press, ISBN 1139462865, 2007.
- 445 [10] J. Fan, T. Hu, Q. Wu, D.-X. Zhou, Consistency analysis of an empirical
446 minimum error entropy algorithm, Applied and Computational Harmonic
447 Analysis 41 (1) (2016) 164–189.
- 448 [11] Z.-C. Guo, D.-H. Xiang, X. Guo, D.-X. Zhou, Thresholded spectral algo-
449 rithms for sparse approximations, Analysis and Applications (2017) 1–23.
- 450 [12] Z.-C. Guo, Y. Ying, D.-X. Zhou, Online regularized learning with pairwise
451 loss functions, Advances in Computational Mathematics (2016) 1–24.

- 452 [13] D. Harrison, D. L. Rubinfeld, Hedonic housing prices and the demand for
453 clean air, *Journal of environmental economics and management* 5 (1) (1978)
454 81–102.
- 455 [14] T. Hu, J. Fan, Q. Wu, D.-X. Zhou, Regularization schemes for minimum
456 error entropy principle, *Analysis and Applications* 13 (04) (2015) 437–455.
- 457 [15] P. Kar, B. Sriperumbudur, P. Jain, H. Karnick, On the Generalization
458 Ability of Online Learning Algorithms for Pairwise Loss Functions, in:
459 Proceedings of The 30th International Conference on Machine Learning,
460 441–449, 2013.
- 461 [16] J. Lin, L. Rosasco, D.-X. Zhou, Iterative regularization for learning with
462 convex loss functions, *Journal of Machine Learning Research* 17 (77) (2016)
463 1–38.
- 464 [17] J. Lin, D.-X. Zhou, Learning theory of randomized Kaczmarz algorithm,
465 *Journal of Machine Learning Research* 16 (2015) 3341–3365.
- 466 [18] J. Lin, D.-X. Zhou, Online Learning Algorithms can Converge Comparably
467 Fast as Batch Learning, *IEEE Transactions on Neural Networks and*
468 *Learning Systems* .
- 469 [19] R. Meir, T. Zhang, Generalization error bounds for Bayesian mixture al-
470 gorithms, *The Journal of Machine Learning Research* 4 (2003) 839–860.
- 471 [20] S. Mukherjee, Q. Wu, Estimation of gradients and coordinate covariation
472 in classification, *The Journal of Machine Learning Research* 7 (2006) 2481–
473 2514.
- 474 [21] W. Rejchel, On ranking and generalization bounds, *The Journal of Machine*
475 *Learning Research* 13 (1) (2012) 1373–1392.
- 476 [22] O. Shamir, T. Zhang, Stochastic Gradient Descent for Non-smooth Opti-
477 mization: Convergence Results and Optimal Averaging Schemes, in: Pro-
478 ceedings of the 30th International Conference on Machine Learning, 71–79,
479 2013.

- 480 [23] L. Shi, Y.-L. Feng, D.-X. Zhou, Concentration estimates for learning with
481 ℓ_1 -regularizer and data dependent hypothesis spaces, *Applied and Compu-*
482 *tational Harmonic Analysis* 31 (2) (2011) 286–302.
- 483 [24] S. Smale, Y. Yao, Online learning algorithms, *Foundations of Computa-*
484 *tional Mathematics* 6 (2) (2006) 145–170.
- 485 [25] I. Steinwart, A. Christmann, *Support Vector Machines*, Springer Science
486 Business Media, ISBN 0387772421, 2008.
- 487 [26] V. N. Vapnik, V. Vapnik, *Statistical learning theory*, vol. 1, Wiley New
488 York, 1998.
- 489 [27] Y. Wang, R. Khardon, D. Pechyony, R. Jones, Generalization Bounds for
490 Online Learning Algorithms with Pairwise Loss Functions., in: *COLT*,
491 vol. 23, 13–1, 2012.
- 492 [28] K. Q. Weinberger, J. Blitzer, L. K. Saul, Distance metric learning for large
493 margin nearest neighbor classification, in: *Advances in neural information*
494 *processing systems*, 1473–1480, 2005.
- 495 [29] Q. Wu, D.-X. Zhou, Learning with sample dependent hypothesis spaces,
496 *Computers & Mathematics with Applications* 56 (11) (2008) 2896–2907.
- 497 [30] Y. Ying, Q. Wu, C. Campbell, Learning the coordinate gradients, *Advances*
498 *in Computational Mathematics* 37 (3) (2012) 355–378.
- 499 [31] Y. Ying, D. X. Zhou, Online regularized classification algorithms, *IEEE*
500 *Transactions on Information Theory* 52 (11) (2006) 4775–4788.
- 501 [32] Y. Ying, D.-X. Zhou, Online pairwise learning algorithms, *Neural compu-*
502 *tation* 28 (4) (2016) 743–777.
- 503 [33] Y. Ying, D.-X. Zhou, Unregularized Online Learning Algorithms with Gen-
504 eral Loss Functions, *Applied and Computational Harmonic Analysis* 42 (2)
505 (2017) 224–244.

506 [34] P. Zhao, R. Jin, T. Yang, S. C. Hoi, Online AUC maximization, in: Pro-
 507 ceedings of the 28th International Conference on Machine Learning (ICML-
 508 11), 233–240, 2011.

509 **Appendix A. Appendix for Proving (3.14)**

First note that

$$\sum_{t=T-k+1}^T t^{-q^*} \leq \int_{T-k}^T x^{-q^*} dx \leq \frac{T^{1-q^*} - (T-k)^{1-q^*}}{1-q^*}, \quad \text{when } q^* \neq 1.$$

When $0 < q^* < 1$, for $k \leq \frac{T}{2}$,

$$\sum_{t=T-k}^T t^{-q^*} \leq (T-k)^{-q^*} (k+1) \leq 2^{q^*} T^{-q^*} (k+1),$$

and for $k > \frac{T}{2}$

$$\sum_{t=T-k}^T t^{-q^*} \leq \frac{T^{1-q^*} - (T-k)^{1-q^*}}{1-q^*} + (T-k)^{-q^*} \leq \frac{T^{1-q^*}}{1-q^*}.$$

510 It thus follows that

$$\begin{aligned} & \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k}^T t^{-q^*} \\ & \leq \sum_{k \leq T/2} \frac{1}{k(k+1)} 2^{q^*} T^{-q^*} (k+1) + \sum_{T-1 \geq k > T/2} \frac{1}{k(k+1)} \frac{T^{1-q^*}}{1-q^*} \\ & \leq \left(2^{q^*+1} + \frac{2}{1-q^*} \right) (\log T) T^{-q^*}. \end{aligned}$$

511 When $q^* = 1$, we have

$$\begin{aligned} \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k}^T t^{-q^*} & \leq \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \frac{k+1}{T-k} = \frac{1}{T} \sum_{k=1}^{T-1} \left\{ \frac{1}{k} + \frac{1}{T-k} \right\} \\ & \leq 4(\log T) T^{-1}. \end{aligned}$$

When $2 > q^* > 1$, for $k \leq \frac{T}{2}$,

$$\sum_{t=T-k}^T t^{-q^*} \leq (T-k)^{-q^*} (k+1) \leq 2^{q^*} T^{-q^*} (k+1),$$

and for $k > \frac{T}{2}$

$$\sum_{t=T-k}^T t^{-q^*} \leq \frac{(T-k)^{1-q^*} - T^{1-q^*}}{q^* - 1} + (T-k)^{-q^*} \leq \frac{q^*}{q^* - 1}.$$

512 Therefore, we have

$$\begin{aligned} & \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k}^T t^{-q^*} \\ & \leq 2^{q^*} T^{-q^*} \sum_{k \leq T/2} \frac{1}{k} + \frac{q^*}{q^* - 1} \sum_{T-1 \geq k > T/2} \frac{1}{k(k+1)} \\ & \leq 2^{q^*+1} T^{-q^*} \log T + \frac{2q^*}{q^* - 1} T^{-1} \\ & \leq \frac{2^{q^*+1} + 2q^*}{q^* - 1} T^{-1}. \end{aligned}$$