

Projection-free Online Learning over Strongly Convex Sets

Yuanyu Wan¹, Lijun Zhang^{1,2,*}

¹Nanjing University of Aeronautics and Astronautics, Nanjing 210023, C
²Peking University, Beijing 100871, G
 {yuanwu, lijunzhang}@pku.edu.cn

Abstract

Online learning over strongly convex sets has been studied extensively in the literature. The best known algorithm is the Online Gradient Descent (OGD) algorithm, which achieves a regret of $O(\sqrt{T})$ over a strongly convex set. However, OGD requires the projection of the subgradients onto the set, which is computationally expensive. In this paper, we propose a new algorithm, the Online Frank-Wolfe (OFW) algorithm, which achieves a regret of $O(\sqrt{T})$ over a strongly convex set without the need for projections. Our algorithm is based on the Frank-Wolfe (FW) algorithm, which is a first-order method for minimizing a convex function over a convex set. We show that the OFW algorithm achieves a regret of $O(\sqrt{T})$ over a strongly convex set, which is the same as the OGD algorithm. Our algorithm is simple and efficient, and it can be applied to a wide range of problems.

Introduction

Online learning over strongly convex sets has been studied extensively in the literature. The best known algorithm is the Online Gradient Descent (OGD) algorithm, which achieves a regret of $O(\sqrt{T})$ over a strongly convex set. However, OGD requires the projection of the subgradients onto the set, which is computationally expensive. In this paper, we propose a new algorithm, the Online Frank-Wolfe (OFW) algorithm, which achieves a regret of $O(\sqrt{T})$ over a strongly convex set without the need for projections. Our algorithm is based on the Frank-Wolfe (FW) algorithm, which is a first-order method for minimizing a convex function over a convex set. We show that the OFW algorithm achieves a regret of $O(\sqrt{T})$ over a strongly convex set, which is the same as the OGD algorithm. Our algorithm is simple and efficient, and it can be applied to a wide range of problems.

$$R(T) = \sum_{t=1}^T f_t(\mathbf{x}_t) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x})$$

Online learning over strongly convex sets has been studied extensively in the literature. The best known algorithm is the Online Gradient Descent (OGD) algorithm, which achieves a regret of $O(\sqrt{T})$ over a strongly convex set. However, OGD requires the projection of the subgradients onto the set, which is computationally expensive. In this paper, we propose a new algorithm, the Online Frank-Wolfe (OFW) algorithm, which achieves a regret of $O(\sqrt{T})$ over a strongly convex set without the need for projections. Our algorithm is based on the Frank-Wolfe (FW) algorithm, which is a first-order method for minimizing a convex function over a convex set. We show that the OFW algorithm achieves a regret of $O(\sqrt{T})$ over a strongly convex set, which is the same as the OGD algorithm. Our algorithm is simple and efficient, and it can be applied to a wide range of problems.

*Lijun Zhang is the corresponding author. Email: lijunzhang@pku.edu.cn

A	E	C	L	E	C	\mathcal{K}	R	B
OFW								$O(T^{3/4})$
LLOO-OCO								$O(\bar{T})$
LLOO-OCO								$O(\log T)$
F OGD								$O(\bar{T})$
F OGD								$O(\log T)$
OSPF								$O(T^{2/3})$
OFW L S ()								$O(T^{2/3})$
SC-OFW ()								$O(T^{2/3})$
SC-OFW ()								$O(\bar{T})$

T 1: C (H 2016), LLOO-OCO (G H 2016), F OGD (L K 2019), OSPF (H M 2020), OFW (H K 2012);

OFW (SC-OFW) (2016) $O(\log T)$ $O(T^{2/3})$ $O(\bar{T})$

Related Work

OFW (H K 2012; H 2016) $O(T^{3/4})$ $O(\bar{T})$ $O(\log T)$ $O(T^{2/3})$

OFW (F M 2005; B 2015), $O(T^{4/5})$ $O(T^{3/4})$ $\tilde{O}(T^{2/3})$ $O(T)$

OFW (2019) $O(T)$ $O(\bar{T})$ $O(T^{3/4})$ $O(T^{2/3})$ $O(\bar{T})$

Main Results

OFW SC-OFW

Preliminaries

$x, y \in \mathcal{K}$

(Bertsekas and Nedic, 2004),

Definition 1 Let $f(\mathbf{x}) : \mathcal{K} \rightarrow \mathbb{R}$ be a function over \mathcal{K} . It is called β -smooth over \mathcal{K} if for all $\mathbf{x}, \mathbf{y} \in \mathcal{K}$

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Definition 2 Let $f(\mathbf{x}) : \mathcal{K} \rightarrow \mathbb{R}$ be a function over \mathcal{K} . It is called α -strongly convex over \mathcal{K} if for all $\mathbf{x}, \mathbf{y} \in \mathcal{K}$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Let $f(\mathbf{x}) : \mathcal{K} \rightarrow \mathbb{R}$ be α -strongly convex and $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$. Then (2015)

$$\frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \quad (1)$$

$$\|\nabla f(\mathbf{x})\|_2 \geq \sqrt{\frac{\alpha}{2} (f(\mathbf{x}) - f(\mathbf{x}^*))}. \quad (2)$$

Let $\mathcal{K} \subseteq \mathbb{E}$ be a convex set. Then (Liu et al., 1979; Gorbunov et al., 2015; Recht and Ben-Ben-Shimon, 2017; Duchi et al., 2019).

Definition 3 A convex set $\mathcal{K} \subseteq \mathbb{E}$ is called α -strongly convex with respect to a norm $\|\cdot\|$ if for any $\mathbf{x}, \mathbf{y} \in \mathcal{K}$, $\gamma \in [0, 1]$ and $\mathbf{z} \in \mathbb{E}$ such that $\|\mathbf{z}\| = 1$, it holds that

$$\gamma \mathbf{x} + (1 - \gamma) \mathbf{y} + \gamma(1 - \gamma) \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2 \mathbf{z} \in \mathcal{K}.$$

Assume $\mathcal{K} \subseteq \mathbb{R}^d$ is α -strongly convex with respect to the ℓ_p norm (2015),

$$\mathcal{K} = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\|_p \leq r\}$$

then \mathcal{K} is α -strongly convex with respect to the ℓ_2 norm (Gorbunov et al., 2015; OCO (Sridharan et al., 2011; Hesterberg et al., 2016)).

Assumption 1 The diameter of the convex decision set \mathcal{K} is bounded by D , i.e.,

$$\|\mathbf{x} - \mathbf{y}\|_2 \leq D$$

for any $\mathbf{x}, \mathbf{y} \in \mathcal{K}$.

Assumption 2 At each round t , the loss function $f_t(\mathbf{x})$ is G -Lipschitz over \mathcal{K} , i.e.,

$$f_t(\mathbf{x}) - f_t(\mathbf{y}) \leq G \|\mathbf{x} - \mathbf{y}\|_2$$

for any $\mathbf{x}, \mathbf{y} \in \mathcal{K}$.

Algorithm 1 OFW

- 1: **Input:** \mathcal{K}, η
 - 2: **Initialization:** $\mathbf{x}_1 \in \mathcal{K}$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: **D** fi $F_t(\mathbf{x}) = \eta \sum_{\tau=1}^t \langle f_\tau(\mathbf{x}_\tau), \mathbf{x} \rangle + \|\mathbf{x} - \mathbf{x}_1\|_2^2$
 - 5: $\mathbf{v}_t \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} \langle F_t(\mathbf{x}_t), \mathbf{x} \rangle$
 - 6: $\sigma_t = \operatorname{argmin}_{\sigma \in [0,1]} (\sigma \langle \mathbf{v}_t - \mathbf{x}_t, \mathbf{v}_t - \mathbf{x}_t \rangle + F_t(\mathbf{x}_t) + \sigma^2 \|\mathbf{v}_t - \mathbf{x}_t\|_2^2)$
 - 7: $\mathbf{x}_{t+1} = \mathbf{x}_t + \sigma_t (\mathbf{v}_t - \mathbf{x}_t)$
 - 8: **end for**
-

OFW with Line Search

Let $\mathcal{K} \subseteq \mathbb{R}^d$ be a convex set. Then (2016)

$$\mathbf{v} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} \langle F_t(\mathbf{x}_t), \mathbf{x} \rangle$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \sigma_t (\mathbf{v} - \mathbf{x}_t)$$

$$F_t(\mathbf{x}) = \eta \sum_{\tau=1}^t \langle f_\tau(\mathbf{x}_\tau), \mathbf{x} \rangle + \|\mathbf{x} - \mathbf{x}_1\|_2^2 \quad (3)$$

Then (2016), OFW with line search (Gorbunov et al., 2015),

$$\sigma_t = \operatorname{argmin}_{\sigma \in [0,1]} (\sigma \langle \mathbf{v}_t - \mathbf{x}_t, \mathbf{v}_t - \mathbf{x}_t \rangle + F_t(\mathbf{x}_t) + \sigma^2 \|\mathbf{v}_t - \mathbf{x}_t\|_2^2)$$

Then (2016), OFW with line search (Gorbunov et al., 2015),

Lemma 1 Let \mathcal{K} be an α_K -strongly convex set with respect to the ℓ_2 norm. Let $\mathbf{x}_t^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} F_{t-1}(\mathbf{x})$ for any $t \in [T + 1]$, where $F_t(\mathbf{x})$ is defined in (3). Then, for any $t \in [T + 1]$, Algorithm 1 with $\eta = \frac{D}{2G(T+2)^{2/3}}$ has

$$F_{t-1}(\mathbf{x}_t) - F_{t-1}(\mathbf{x}_t^*) \leq \epsilon_t = \frac{C}{(t+2)^{2/3}}$$

where $C = \max\left(4D^2, \frac{4096}{3\alpha_K^2}\right)$.

Then (2016), OFW with line search (Gorbunov et al., 2016, L.3).

$$\begin{aligned}
& \sum_{t=1}^T f_t(\mathbf{x}_t) - \sum_{t=1}^T f_t(\mathbf{x}^*) \\
& \leq \sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle \\
& = \underbrace{\sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_t^* \rangle}_{:=A} + \underbrace{\sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_t^* - \mathbf{x}^* \rangle}_{:=B}.
\end{aligned} \tag{6}$$

Lemma 3 (Lemma 6.6 of Garber and Hazan (2016)) Let $\{f_t(\mathbf{x})\}_{t=1}^T$ be a sequence of loss functions and let $\mathbf{x}_t^* \in \arg\min_{\mathbf{x} \in \mathcal{K}} \sum_{\tau=1}^t f_\tau(\mathbf{x})$ for any $t \in [T]$. Then, it holds that

$$\begin{aligned}
A & = \sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_t^* \rangle \\
& \leq \sum_{t=1}^T \|f_t(\mathbf{x}_t)\|_2 \|\mathbf{x}_t - \mathbf{x}_t^*\|_2 \\
& \leq \sum_{t=1}^T G \sqrt{F_{t-1}(\mathbf{x}_t) - F_{t-1}(\mathbf{x}_t^*)} \\
& \leq \sum_{t=1}^T \frac{G \bar{C}}{(t+2)^{1/3}} \leq \frac{3G \bar{C}(T+2)^{2/3}}{2}.
\end{aligned} \tag{7}$$

$$\sum_{t=1}^T (t+2)^{-1/3} \leq 3(T+2)^{2/3}/2.$$

Lemma 3 (Lemma 6.6 of Garber and Hazan (2016)) Let $\{f_t(\mathbf{x})\}_{t=1}^T$ be a sequence of loss functions and let $\mathbf{x}_t^* \in \arg\min_{\mathbf{x} \in \mathcal{K}} \sum_{\tau=1}^t f_\tau(\mathbf{x})$ for any $t \in [T]$. Then, it holds that

$$\sum_{t=1}^T f_t(\mathbf{x}_t^*) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}) \leq 0.$$

Lemma 3 (Lemma 6.6 of Garber and Hazan (2016)) Let $\{f_t(\mathbf{x})\}_{t=1}^T$ be a sequence of loss functions and let $\mathbf{x}_t^* \in \arg\min_{\mathbf{x} \in \mathcal{K}} \sum_{\tau=1}^t f_\tau(\mathbf{x})$ for any $t \in [T]$. Then, it holds that

$$\begin{aligned}
& \sum_{t=1}^T \tilde{f}_t(\mathbf{x}_{t+1}^*) - \sum_{t=1}^T \tilde{f}_t(\mathbf{x}^*) \leq 0 \\
& \sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_{t+1}^* - \mathbf{x}^* \rangle \\
& \leq (\|\mathbf{x}^* - \mathbf{x}_1\|_2^2 - \|\mathbf{x}_2^* - \mathbf{x}_1\|_2^2) / \eta \\
& \leq D^2 / \eta.
\end{aligned} \tag{8}$$

$$\|\mathbf{x}_2^* - \mathbf{x}_1\|_2^2 \geq 0$$

$$\begin{aligned}
& \|\mathbf{x}_t^* - \mathbf{x}_{t+1}^*\|_2^2 \\
& \leq F_t(\mathbf{x}_t^*) - F_t(\mathbf{x}_{t+1}^*) \\
& = F_{t-1}(\mathbf{x}_t^*) - F_{t-1}(\mathbf{x}_{t+1}^*) + \eta \langle f_t(\mathbf{x}_t), \mathbf{x}_t^* - \mathbf{x}_{t+1}^* \rangle \\
& \leq \eta \|f_t(\mathbf{x}_t)\|_2 \|\mathbf{x}_t^* - \mathbf{x}_{t+1}^*\|_2 \\
& \|\mathbf{x}_t^* - \mathbf{x}_{t+1}^*\|_2 \leq \eta \|f_t(\mathbf{x}_t)\|_2 \leq \eta G.
\end{aligned} \tag{9}$$

$$\eta = \frac{D}{2G(T+2)^{2/3}},$$

$$\begin{aligned}
B & = \sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_t^* - \mathbf{x}^* \rangle \\
& = \sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_{t+1}^* - \mathbf{x}^* \rangle \\
& \quad + \sum_{t=1}^T \langle f_t(\mathbf{x}_t), \mathbf{x}_t^* - \mathbf{x}_{t+1}^* \rangle \\
& \leq \frac{D^2}{\eta} + \sum_{t=1}^T \|f_t(\mathbf{x}_t)\|_2 \|\mathbf{x}_t^* - \mathbf{x}_{t+1}^*\|_2 \\
& \leq \frac{D^2}{\eta} + \eta T G^2 \\
& \leq 2DG(T+2)^{2/3} + \frac{DG(T+2)^{1/3}}{2} \\
& \leq G \bar{C}(T+2)^{2/3} + \frac{G \bar{C}(T+2)^{2/3}}{4}
\end{aligned} \tag{10}$$

$$D \leq \bar{C}/2 \quad (T+2)^{1/3} \geq (T+2)^{2/3} \quad T \geq \frac{D}{\bar{C}}.$$

Proof of Theorem 2

Let $\tilde{f}_t(\mathbf{x}) = \langle f_t(\mathbf{x}_t), \mathbf{x} \rangle + \frac{\lambda}{2} \|\mathbf{x} - \mathbf{x}_t\|_2^2$, $t \in [T]$
 $\mathbf{x}_t^* = \arg\min_{\mathbf{x} \in \mathcal{K}} F_{t-1}(\mathbf{x})$, $t = 2, \dots, T+1$.

$$\begin{aligned}
& \sum_{t=1}^T f_t(\mathbf{x}_t) - \sum_{t=1}^T f_t(\mathbf{x}^*) \\
& \leq \sum_{t=1}^T \left(\langle f_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle - \frac{\lambda}{2} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 \right) \\
& = \sum_{t=1}^T (\tilde{f}_t(\mathbf{x}_t) - \tilde{f}_t(\mathbf{x}^*)) \\
& = \underbrace{\sum_{t=1}^T (\tilde{f}_t(\mathbf{x}_t) - \tilde{f}_t(\mathbf{x}_{t+1}^*))}_{:=A} + \underbrace{\sum_{t=1}^T (\tilde{f}_t(\mathbf{x}_{t+1}^*) - \tilde{f}_t(\mathbf{x}^*))}_{:=B}.
\end{aligned}$$

Lemma 3 (Lemma 6.6 of Garber and Hazan (2016)) Let $\{f_t(\mathbf{x})\}_{t=1}^T$ be a sequence of loss functions and let $\mathbf{x}_t^* \in \arg\min_{\mathbf{x} \in \mathcal{K}} \sum_{\tau=1}^t f_\tau(\mathbf{x})$ for any $t \in [T]$. Then, it holds that

Lemma 4 (Lemma 6.7 of Garber and Hazan (2016)) For any $t \in [T]$, the function $\tilde{f}_t(\mathbf{x}) = \langle f_t(\mathbf{x}_t), \mathbf{x} \rangle + \frac{\lambda}{2} \|\mathbf{x} - \mathbf{x}_t\|_2^2$ is $(G + \lambda D)$ -Lipschitz over \mathcal{K} .

$$\text{B} \quad \text{L} \quad 4, \mathbf{i} \quad t = 3, \dots, T + 1,$$

$$\begin{aligned} & F_{t-1}(\mathbf{x}_{t-1}^*) - F_{t-1}(\mathbf{x}_t^*) \\ &= F_{t-2}(\mathbf{x}_{t-1}^*) - F_{t-2}(\mathbf{x}_t^*) + \tilde{f}_{t-1}(\mathbf{x}_{t-1}^*) - \tilde{f}_{t-1}(\mathbf{x}_t^*) \\ &\leq (G + \lambda D) \|\mathbf{x}_{t-1}^* - \mathbf{x}_t^*\|_2. \end{aligned}$$

$$\text{M} \quad , \quad \sim \quad \sim \quad F_t(\mathbf{x}) \quad t\lambda \quad \sim \quad , \quad \mathbf{i} \\ t = 3, \dots, T + 1,$$

$$\begin{aligned} \|\mathbf{x}_{t-1}^* - \mathbf{x}_t^*\|_2^2 &\leq \frac{2(F_{t-1}(\mathbf{x}_{t-1}^*) - F_{t-1}(\mathbf{x}_t^*))}{(t-1)\lambda} \\ &\leq \frac{2(G + \lambda D) \|\mathbf{x}_{t-1}^* - \mathbf{x}_t^*\|_2}{(t-1)\lambda}. \end{aligned}$$

$$\text{T} \quad \mathbf{i} \quad , \quad \mathbf{i} \quad t = 3, \dots, T + 1,$$

$$\|\mathbf{x}_{t-1}^* - \mathbf{x}_t^*\|_2 \leq \frac{2(G + \lambda D)}{(t-1)\lambda}. \quad (11)$$

$$\text{B} \quad \text{L} \quad 2 \quad 4,$$

$$\begin{aligned} & \sum_{t=2}^T (\tilde{f}_t(\mathbf{x}_t) - \tilde{f}_t(\mathbf{x}_{t+1}^*)) \\ &\leq \sum_{t=2}^T (G + \lambda D) \|\mathbf{x}_t - \mathbf{x}_{t+1}^*\|_2 \\ &\leq (G + \lambda D) \sum_{t=2}^T \|\mathbf{x}_t - \mathbf{x}_t^*\|_2 \\ &\quad + (G + \lambda D) \sum_{t=2}^T \|\mathbf{x}_t^* - \mathbf{x}_{t+1}^*\|_2 \\ &\leq (G + \lambda D) \sum_{t=2}^T \sqrt{\frac{2(F_{t-1}(\mathbf{x}_t) - F_{t-1}(\mathbf{x}_t^*))}{(t-1)\lambda}} \\ &\quad + (G + \lambda D) \sum_{t=2}^T \frac{2(G + \lambda D)}{t\lambda} \\ &\leq (G + \lambda D) \sum_{t=2}^T \sqrt{\frac{2C}{t}} \end{aligned}$$

- Bertsekas, S.; Uppala, V. 2004. *Convex Optimization*. Cambridge University Press.
- Bertsekas, S.; Duchi, O.; Karmali, T.; Parrilo, P. 2015. Bounding the variance of the stochastic gradient. *Proceedings of the 28th Conference on Learning Theory*, 266-278.
- Cohen, J.; Eluder, T.; Li, Q.; Shalunov, L.; Song, C. 2016. Online learning with linearly growing regret. *Proceedings of the 32nd Conference on Uncertainty in Artificial Intelligence*, 122-131.
- Cohen, L.; Eluder, M.; Karmali, A. 2019. Pseudoregularity. *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, 2047-2056.
- Dennis, V. F.; Moré, A. M. 1970. *Approximate Methods in Optimization Problems*. Elsevier.
- Dimitrakopoulos, J. C.; Athanassios, A.; Wang, M. J. 2011. Distributed control of a multi-agent system. *IEEE Transactions on Automatic Control* 57(3): 592-606.
- Dimitrakopoulos, J. C. 1979. Random search optimization. *SIAM Journal on Control and Optimization* 17(2): 187-211.
- Feldman, A. D.; Karagulyan, A. T.; Matas, H. B. 2005. Online learning with linearly growing regret. *Proceedings of the 16th Annual ACM-SIAM Symposium on Discrete Algorithms*, 385-394.
- Ford, M.; Whittle, P. 1956. A stochastic control problem. *Naval Research Logistics Quarterly* 3(1-2): 95-110.
- Gilmer, D.; Hoyer, E. 2015. Fenchel-epigraph. *Proceedings of the 32nd International Conference on Machine Learning*, 541-549.
- Gilmer, D.; Hoyer, E. 2016. A unified view of matrix factorization. *SIAM Journal on Optimization* 26(3): 1493-1528.
- Gilmer, D.; Karmali, B. 2020. IRL. *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, 2196-2206.
- Gilmer, D.; Karmali, B. 2020. R. *ArXiv e-prints* [arXiv:2010.07572](https://arxiv.org/abs/2010.07572).
- Hoyer, E. 2008. S. *Latin American Symposium on Theoretical Informatics*, 306-316.
- Hoyer, E. 2016. I. *Foundations and Trends in Optimization* 2(3-4): 157-325.
- Hoyer, E.; Amari, A.; Karmali, S. 2007. L. *Machine Learning* 69(2): 169-192.
- Hoyer, E.; Karmali, S. 2012. P. *Proceedings of the 29th International Conference on Machine Learning*, 1843-1850.
- Hoyer, E.; Li, L.; Hoyer, H. 2016. V. *Proceedings of the 33rd International Conference on Machine Learning*, 1263-1271.
- Hoyer, E.; Moré, E. 2020. F. *Proceedings of the 33rd Annual Conference on Learning Theory*, 1877-1893.
- Hoyer, S.; Cohen, A.; Moré, M. 2013. O. *52nd IEEE Conference on Decision and Control*, 1484-1489.
- Jain, M. 2013. R. *Proceedings of the 30th International Conference on Machine Learning*, 427-435.
- Lavrov, E. S.; Pukhov, B. T. 1966. C. *USSR Computational mathematics and mathematical physics* 6: 1-50.
- Lavrov, K.; Karmali, A. 2019. P. *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, 1458-1466.
- Lavrov, H.; Wang, C.; Karmali, K. 2018. E. *Advances in Neural Information Processing Systems* 31, 8235-8245.
- Rubinfeld, J.; Wang, J.-K.; Moré, B. 2019. R. *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, 1576-1583.
- Srebro, S. 2007. *Online Learning: Theory, Algorithms, and Applications*. Cambridge University Press.
- Srebro, S. 2011. O. *Foundations and Trends in Machine Learning* 4(2): 107-194.
- Srebro, S.; Shalunov, L. 2007. A. *Machine Learning* 69(2-3): 115-142.
- Wang, J.; Tang, W.-W.; Li, L. 2020. P. *Proceedings of the 37th International Conference on Machine Learning*, 9818-9828.
- Wang, J.; Li, L. 2020. P. *ArXiv e-prints* [arXiv:2010.08177](https://arxiv.org/abs/2010.08177).
- Wang, J.; Li, Q.; Shalunov, L. 2017. A. *Proceedings of the 34th International Conference on Machine Learning*, 3901-3910.

[1] L. J. R. H., 2013. $O(\log T)$ -
 Proceedings of the 30th International
 Conference on Machine Learning, 1121-1129.
 [2] W. P. W. H., S. C. H.; T. 2017. P
 Proceedings of the 34th International Conference
 on Machine Learning, 4054-4062.
 [3] M. 2003. O
 Proceedings of the
 20th International Conference on Machine Learning, 928-936.